

# Terminology Service Development at the Mayo Clinic

**Harold Solbrig**  
Technical Specialist  
Medical Informatics Research  
Mayo Clinic  
Rochester, MN 55901

# Outline

- **Background** – why terminology services?
- **Where we've been** – the evolution of terminology services at the Mayo Clinic
- **Where we are now** – current tools and approaches to terminology services
- **Where we are going** – technologies, distribution, future applications



# Background

# Terminology

**The lexicon of a special subject language reflects the organisational characteristics of the discipline by tending to provide as many lexical units as there are concepts...**

**Juan C. Sager, *A Practical Course in Terminology Processing*. John Benjamins, 1990**

# Terminology

**The items which are characterised by special reference within a discipline are the ‘terms’ of the discipline, and collectively they form its ‘terminology’; those which function in general reference over a variety of sublanguages are simply called ‘words’, and their totality the ‘vocabulary’.**

**Juan C. Sager. *A Practical Course in Terminology Processing***

# Evolutionary Steps of a Terminology

- 1) **Everyday words (vocabulary)** - the differentiating knowledge of the trade or profession is in the formative stage.
- 2) **Overloaded words and/or acronyms**  
- As the knowledge increases, it becomes cumbersome to continue to use full phrases to express concepts. Common phrases are shortened:  
*“kernel, heap, garbage collection, SCRAM, LASER, ...”*

# Evolutionary Steps of a Terminology

- 3) **Formalized nomenclature** - as the need for precision and detail increases, the management of haphazard wording becomes prohibited. When practical, formalized naming rules are established. *Linnaeus system, chemical names, SNOMED, etc.*
- 4) **Coding and classification schemes** - as specialties and different views emerge, the need to classify, categorize and cross-reference becomes important. *ICD, etc.*

## Evolutionary Steps of a Terminology

- 5) **Thesauri** - Boundaries between specialties change and cross references become necessary between terminologies of different specialties. *UMLS*
- 6) **Reference Terminologies** - Thesauri become unwieldy and too imprecise. A new, synthetic, atomic conceptual organization is formed as a reference point and focus. *SNOMED-CT, Read Codes III*



# Terminological Categories

## **“First Generation” terminological systems**

- Typically targeted for paper based information systems (not “IT-enabled”)
- Simple hierarchies and organization
- Expensive to maintain and reuse

Angelo Rossi Mori, et. al. *Standards to support development of terminological systems for healthcare telematics.*

# Terminological Categories

## “Second Generation” systems

- Compositional systems, based on a *categorical structure* (semantic categories, semantic links [associations between categories])
- Dynamic organization, systematic description of a subject field
- Flexible and extensible
- Limited semantic based processing

# Terminological Categories

- **“Third Generation” systems**
  - Based on a formal model (a set of symbols and a set of formal rules)
  - Dynamic w/ multiple hierarchies
  - Content and updates are formally validated
  - Complete semantic based processing (behavior is independent of names)

# What is a terminology?

## Key characteristics:

- Set of terms, definitions and relationships for a (relatively) non-ambiguous partitioning of the conceptual space of a specialized subject area or discipline.
- NOT necessarily related to computerized data processing (or even data processing, period)
- A formal shared context for communication among members of a specialty or trade.

# What is a terminology?

- **Key characteristics**
  - Term  $\longleftrightarrow$  concept mapping
  - Additional entry phrases including
    - Lexical variants
    - Synonyms
    - Similar or related phrases
  - Intrinsic definitions, annotations, etc.
  - Extrinsic definitions in form of taxonomy / ontology / semantic net

# Terminology vs. Ontology

- **The word “Ontology”, as it is used today refers to the DL-based organization of ‘concepts’**
- **Focus is on the formal organization**
- **Lexical/linguistic section is underspecified**
  - **Attributed definitions**
  - **Terms in multi-languages and contexts**
- **Behavioral characteristics are strictly DL – no rules on how to find a node given an input string...**



# Uses of Terminology in Clinical Practice

# Terminology in Clinical Practice

**Code Sets** – Lists of codes used to fill out forms, data entry, etc.

- Drawn from small to medium size lists
- Typically local to institution
- May not cross databases or applications



# Code Sets in Forms

10. Complete the following information if the isolate is *vibrio cholerae* 01 or 0139:

## Serotype (452) (check one)

- ☐ Inaba (1)
- ☐ Ogawa (2)
- ☐ Hikojima (3)
- ☐ Not Done (4)
- ☐ Unk. (9)

## Biotype (check one)

- ☐ El Tor (1)
- ☐ Classical (2)
- ☐ Not Done (3)
- ☐ Unk. (9)

Patient home state: \_ \_

## 4. Sex: (68)

- ☐ M (1)
- ☐ F (2)
- ☐ Unk. (9)

## The patient has been enrolled at:

- |   |               |
|---|---------------|
| 1 | NIH-sponsored |
| 2 | Other         |
| 3 | None          |
| 9 | Unknown       |

# Terminology in Clinical Practice

**Classifications** - Codes used to summarize information for the purpose of QA, reporting, reimbursement, etc.

- ICD-9-CM, CPT4, ...
- (Usually) Redundant information
- In use since (at least) the 16<sup>th</sup> century  
London Bills of Mortality

# Terminology in Clinical Practice

## Classifications

- **Shape and size of “buckets” depends on intended use**
  - “Killed by several Accidents”
  - “King’s Evil”
  - “Frightened”
  - “Crushed by falling aircraft in terrorism [attack]”
  - “concussion with more than 24 hours loss of consciousness and return to pre-existing conscious level”

# Terminology in Clinical Practice

**Indices** – Codes used to summarize the content of medical records for the purpose of research and retrieval

- MESH, ICD-9-CM, HICDA, ...

# Terminology in Clinical Practice

**Metadata** – Codes that describe the format and content of databases, files, forms, etc.

- Enables sharing of information across institutions
- Enables sharing of information within institutions across time
- Only recently becoming formalized

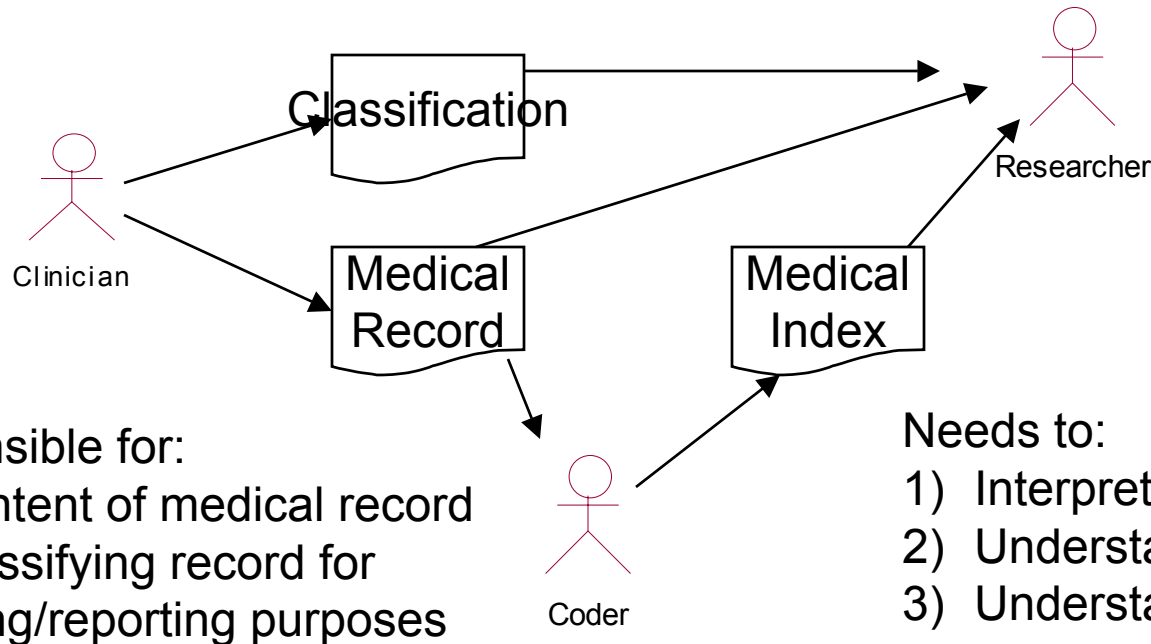
# Terminology in Clinical Practice

- **Code Sets**
  - Often restrictive, incomplete
  - Lack of compositional structure
  - Not applicable in many settings (free text, quantitative data, etc.)
- **Classifications**
  - Granularity depends on the context
  - Rarely matches the level of specificity needed to accurately record clinical information

# Terminology in Clinical Practice

- **Indices**
  - Classification after the fact
  - Can be labor intensive
  - Balance must be maintained between granularity and cost
  - Cannot anticipate unexpected requests
    - AIDS symptoms
    - Terrorism related events

# Terminology in Clinical Practice



Responsible for:

- 1) Content of medical record
- 2) Classifying record for billing/reporting purposes

Needs to:

- 1) Interpret record contents
- 2) Understand classification process
- 3) Understand indexing process

Responsible for:

- 1) Interpreting medical record
- 2) Classifying record for research purposes

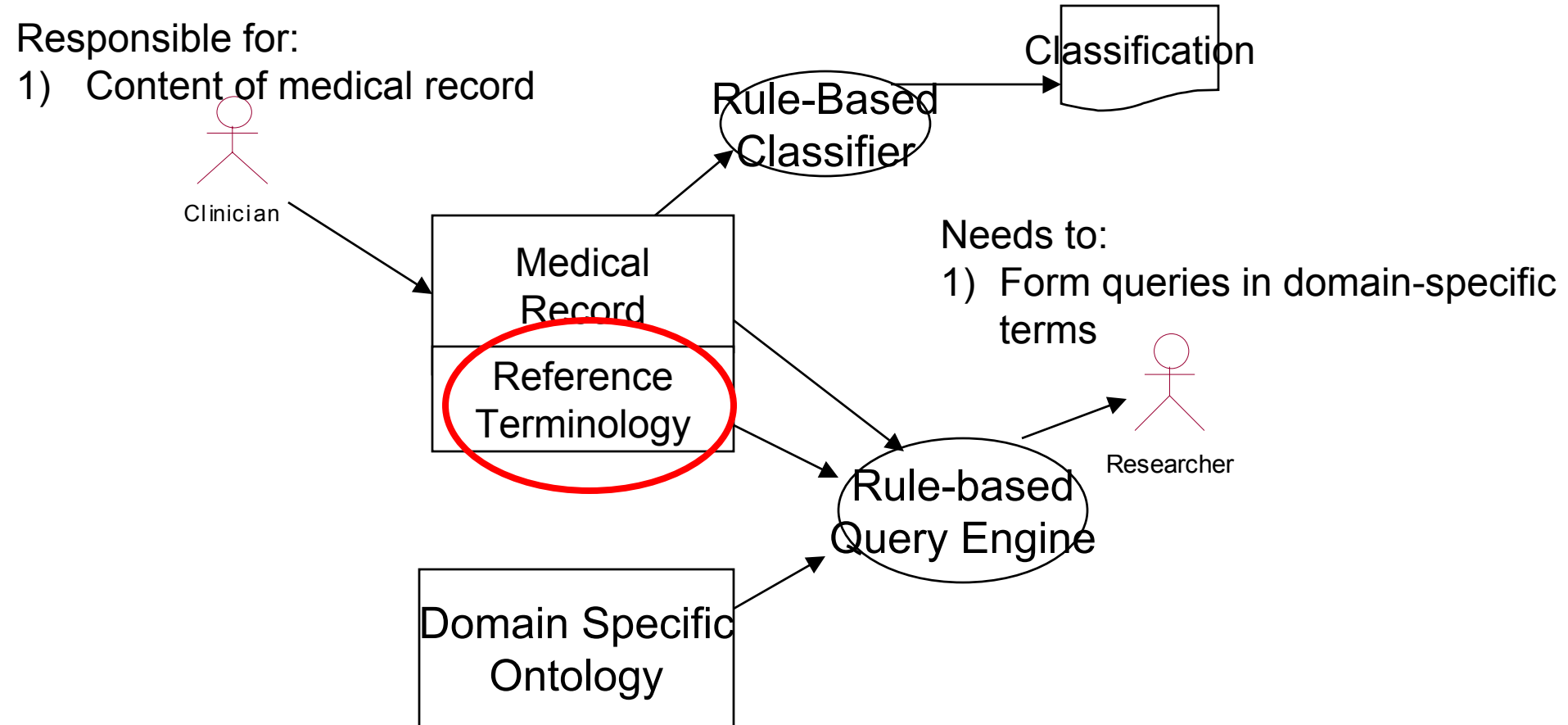




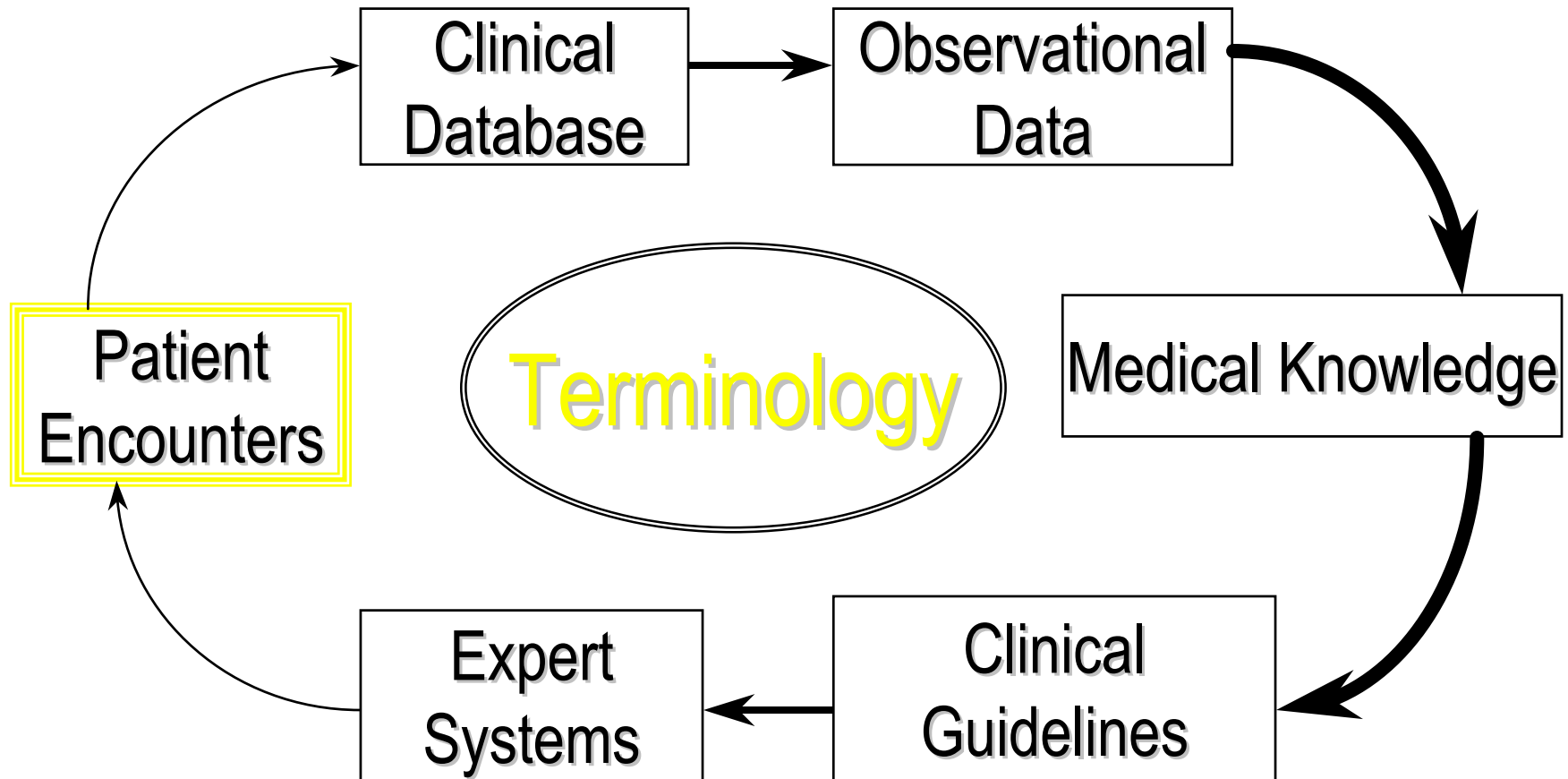
## Some ways that the process can fail

- Incomplete or inconsistent knowledge of classification rules
- Clinician resource time is scarce – billing record is often perceived as a part of the clinical record
- Indexing is resource intensive
- Indexing process depends on what is known at the time
- Researcher has to have intimate understanding of all parts of the process

# Reference Terminology



# Heritage of Continuous Improvement



# Reference Terminology

- **Must represent fine level of clinical detail**
- **Coverage must be broad enough to span an entire discipline**
- **Must be well defined**
- **Must be compositional in nature**
  - **Post-coordination rules**
  - **Rules for determine compositional equivalence**



# Characteristics of a Reference Terminology

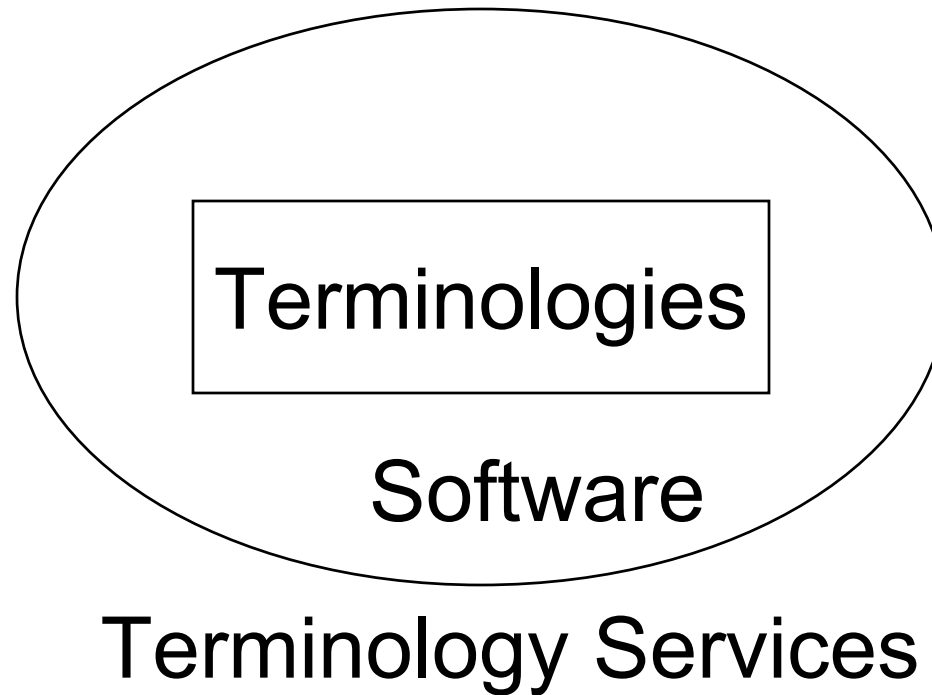
- Requires computers and software to be effectively managed (3<sup>rd</sup> generation system)
- Terminology = content + software

# Requirements for Specifying Software Behavior

- 1. Information model** – What are the entities that are manipulated and how are they related?
- 2. Requirements model** – What questions will the software need answer?
- 3. Behavioral model** – How do request various behaviors, what do they do and how do they respond?



# Reference Terminology





# Where we've been



# Healthcare Data Dictionary (Stan Huff / 3M)

- Evolved from the HELP PTXT system
- Software included:
  - Oracle DB running on Unix
  - Tuxedo transaction management system
  - OLE / C++ based client object system
- Reasonably successful
- Proprietary – content and algorithms didn't port

# Lexicon Query Services (LQS)

- **Developed under auspices of the Object Management Group (OMG)**
- **Read-only – no authoring**
- **Included:**
  - **Information Model of Terminology**
  - **Behavioral Model**
  - **Implementation Specification targeted for the Common Object Request Broker Architecture (CORBA)**
    - **Syntax specification in IDL**
    - **Architecture required significant changes in object/attribute layout**

# LQS Specification

- Schema for globally unique identifiers
- Validate concept code
- Lookup Concept codes by word(s) / string / pattern
- Lookup concept text for a given context / language / lexical type
- Lookup definitions / comments / instructions
- List concepts that have a specified relationship with a supplied concept
- Determine whether two concepts are related
- Reduce a composite expression to a canonical form
- Compare two composite expressions
- ...

# LQS

- Published in 1998
- Not widely used or adopted (outside of 3M)
  - Perceived (and actual) complexity of the specification
  - Not easy to implement
    - Services were not trivial to implement
    - No reference implementation
    - CORBA was difficult and expensive to work with
- Model and functional requirements are still reasonably definitive

# Mayo Terminology Services (MTS)

## (Chris Chute / Mayo)

- Extension of Lexicon Query Services
- Purpose was to provide a complete “breadboard” of terminology components
- Added lexical/linguistic capabilities
  - Spell correction (word locator)
  - Word stemming (using LVG)
  - Word and Phrase completion
  - Plesionymy (words and phrases that could have a very similar meaning in a given context)
  - Candidate term matching

# MTS 2000 Implementation

- Java based
- Used JDBC back end (Oracle / Sybase)
- Used SNOMED-RT Database model
- Multiple Implementation Architectures
  - CORBA
    - COM/DCOM bridge
    - Perl Bridge
    - Python Bridge
  - RMI
  - Straight Java Objects

## MTS - Lessons Learned

Usefulness within Clinic depended on adoption by internal and external software providers

- Y2K came first
- Waited for HIPAA regulations to solidify
- Needed today's external incentives
  - HIPAA
  - Terrorism Surveillance
  - Bioinformatics



# MTS Lessons Learned

Usefulness within Mayo Clinic depended on adoption by internal and external software providers

- Specification needs to be
  - Standards based
  - (Relatively) easy to implement





# Where we are today

## Where we are today

- **Open Terminology Services**
- **HL7 CTS Specification**
- **HL7 Terminology Tools**
- **LDAP Back End**
- **OTS Using LDAP and Lucene**
- **NLP processing of medical records and term source**

# Open Terminology Services

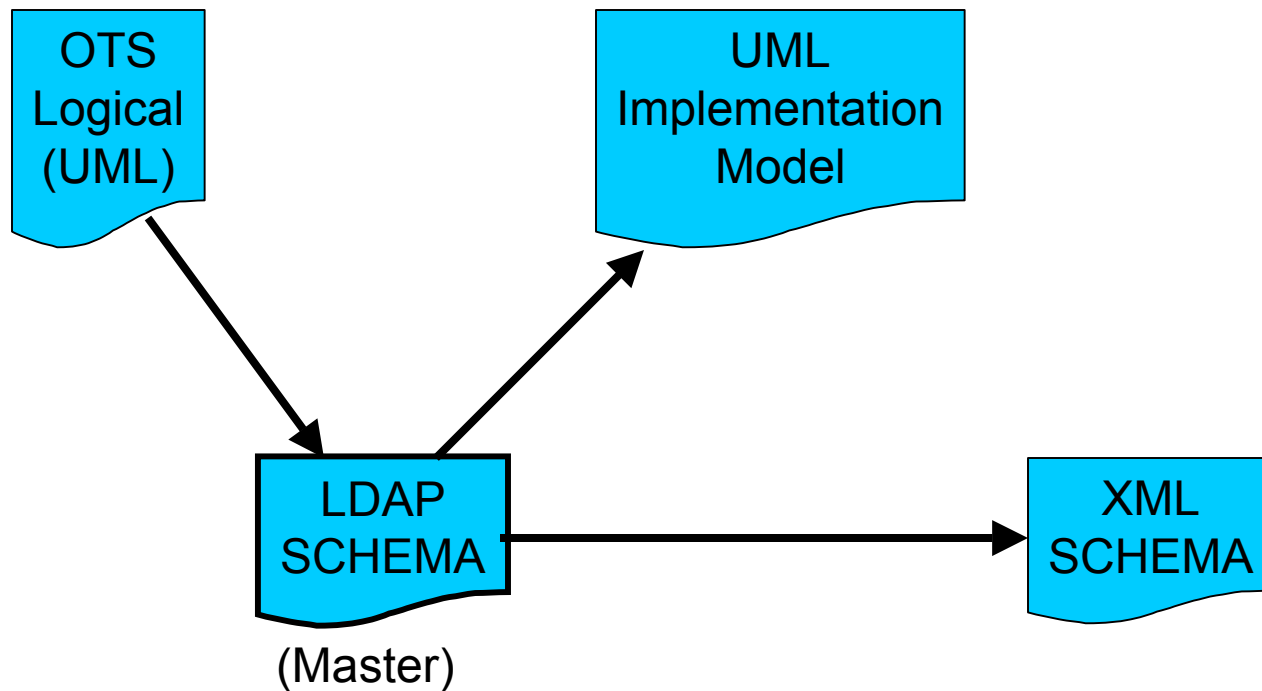
- Refactoring of MTS
- Standards Based:
  - CTS through HL7
  - “Open Source” approach
    - Java Reference implementation
    - LDAP back end
    - ... serious clash w/ Mayo culture

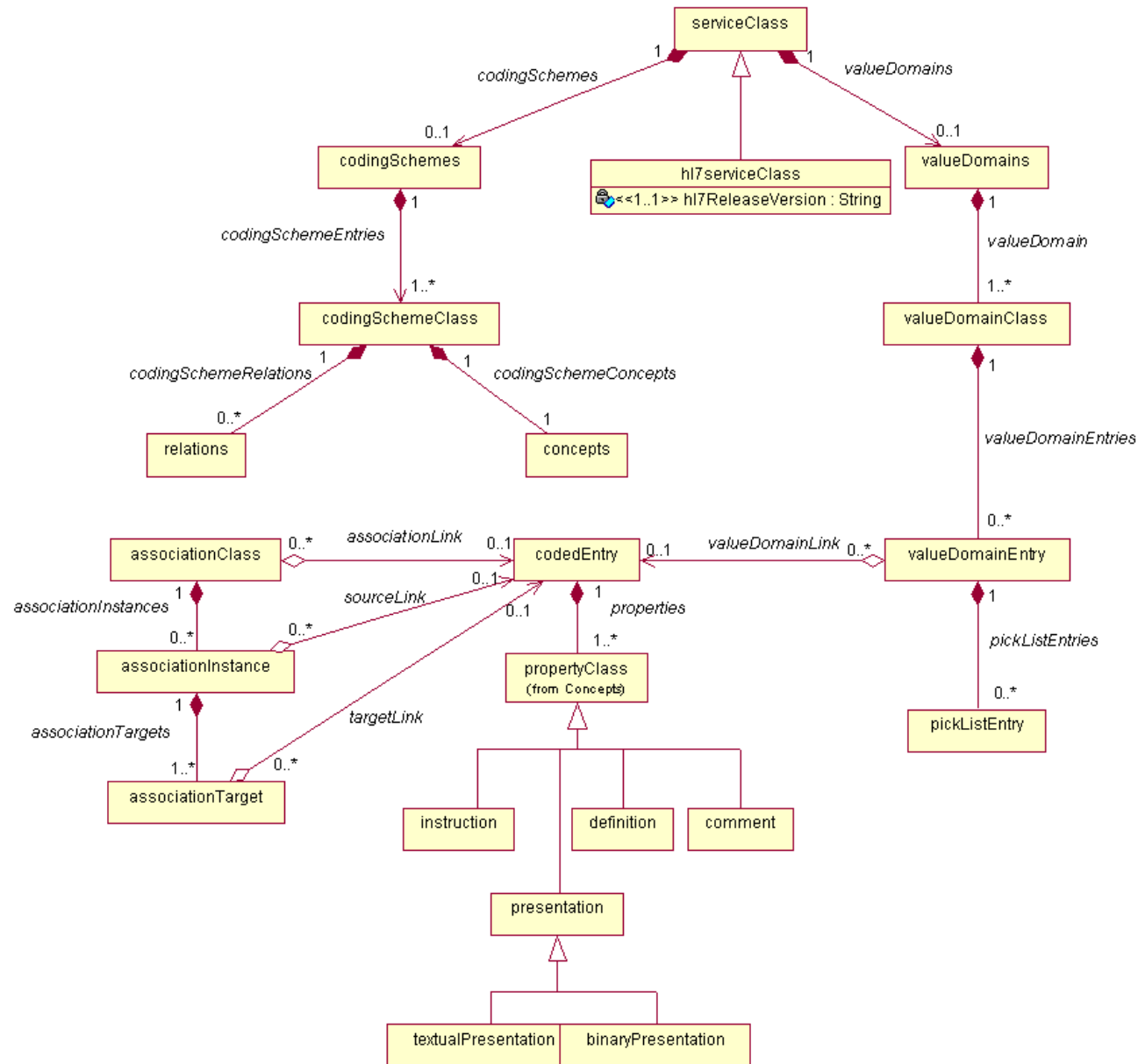
# OTS Components

- Model
  - Abstract ('logical') model of information content
  - Implementation models – one per technology
- Content
  - Terminology content deployed in various technologies
- Software
  - Browsing and Implementation Tools
  - Distribution and deployment tools
  - Editing and revision tools



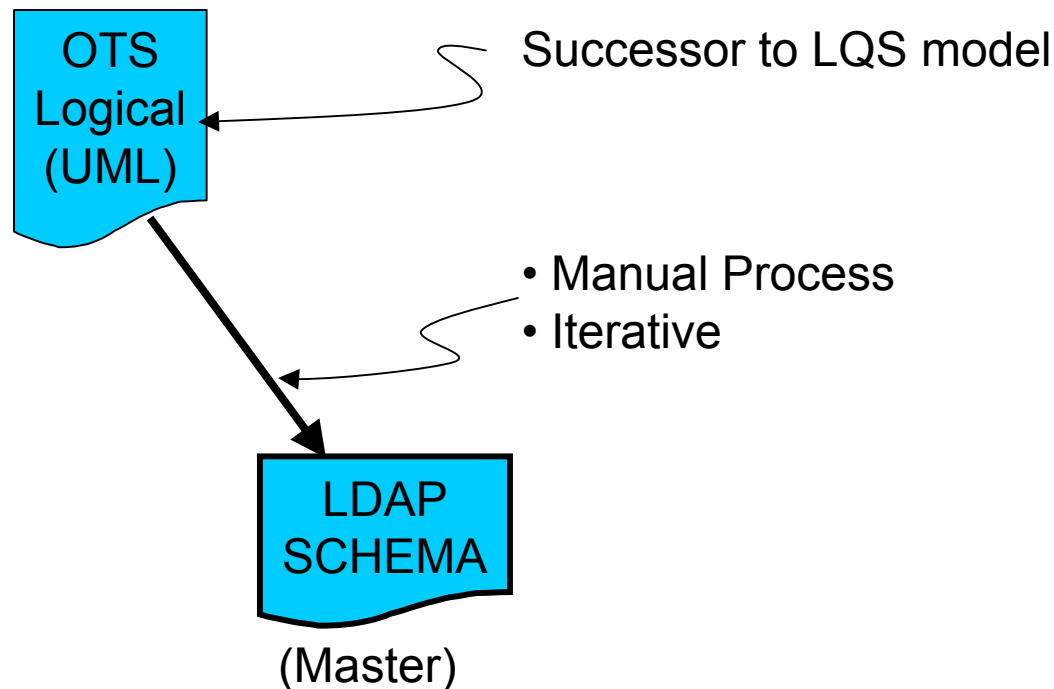
# OTS Model



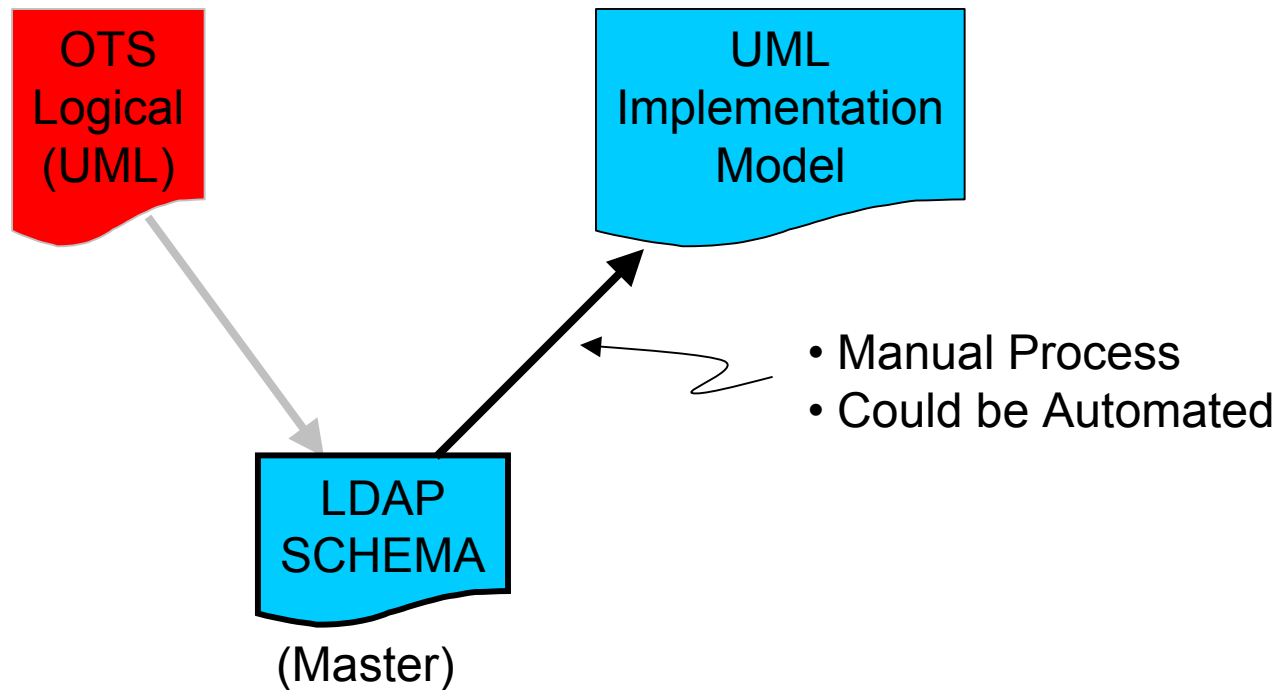


# OTS Model

## Logical Model & Schema



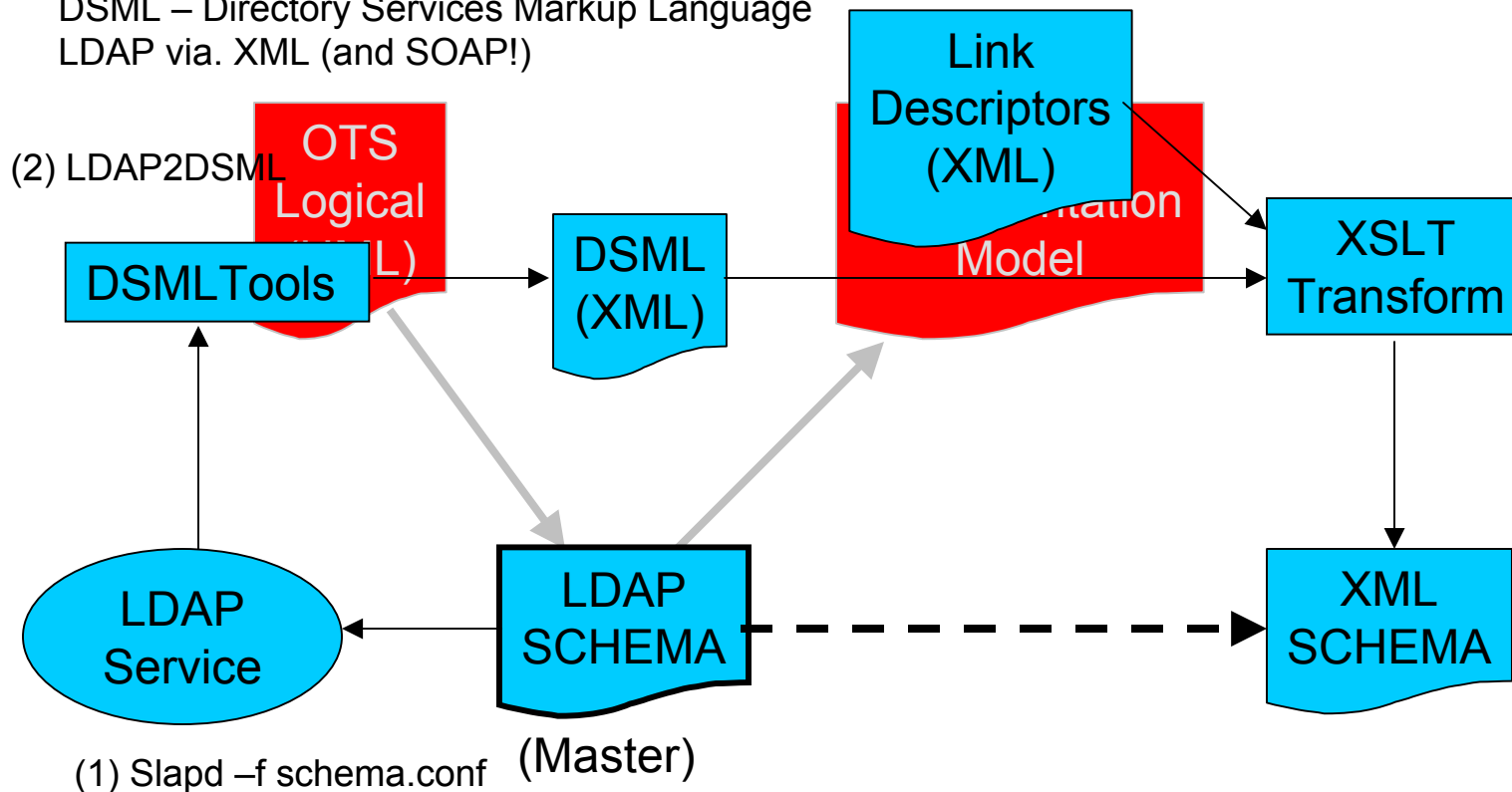
# OTS Model Implementation Model

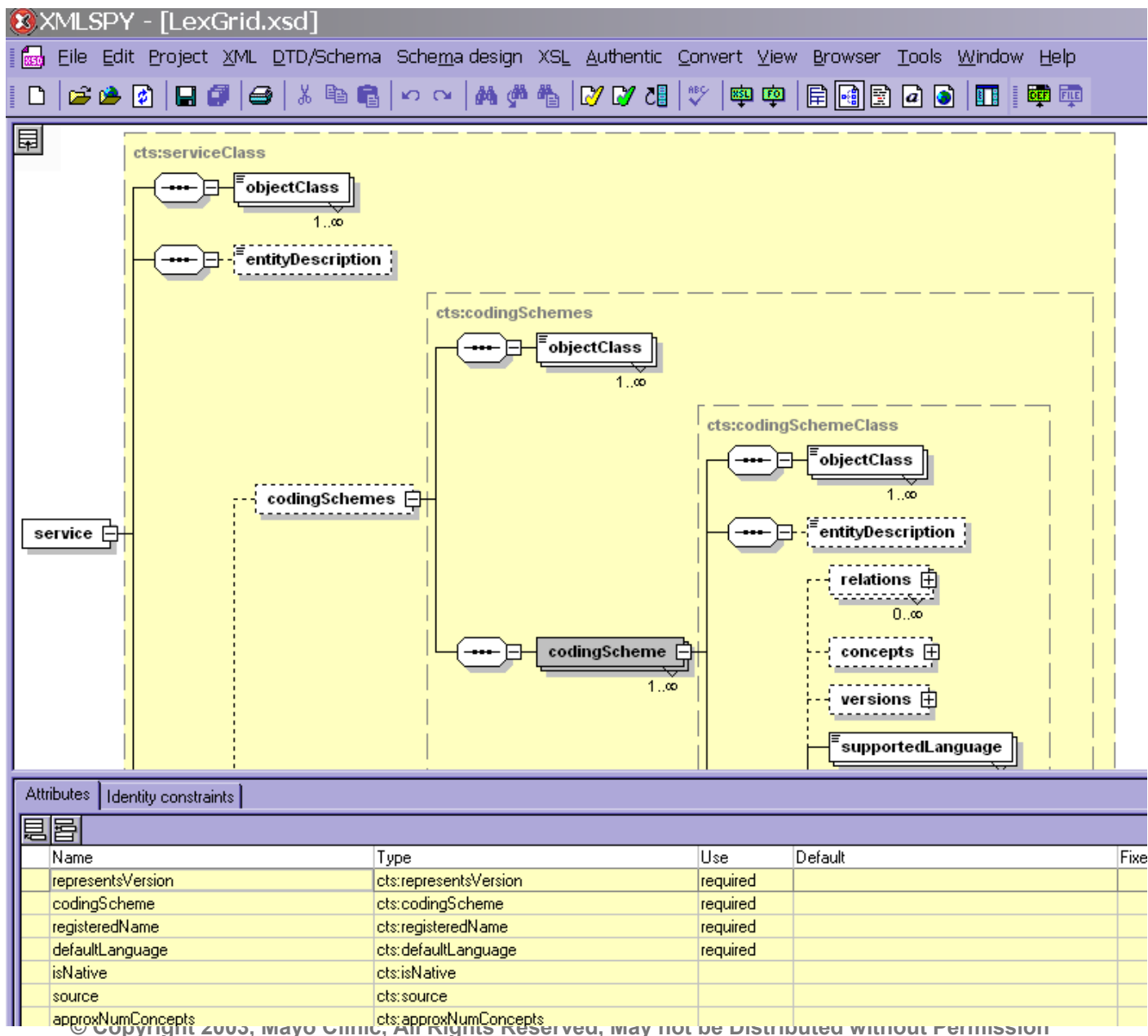




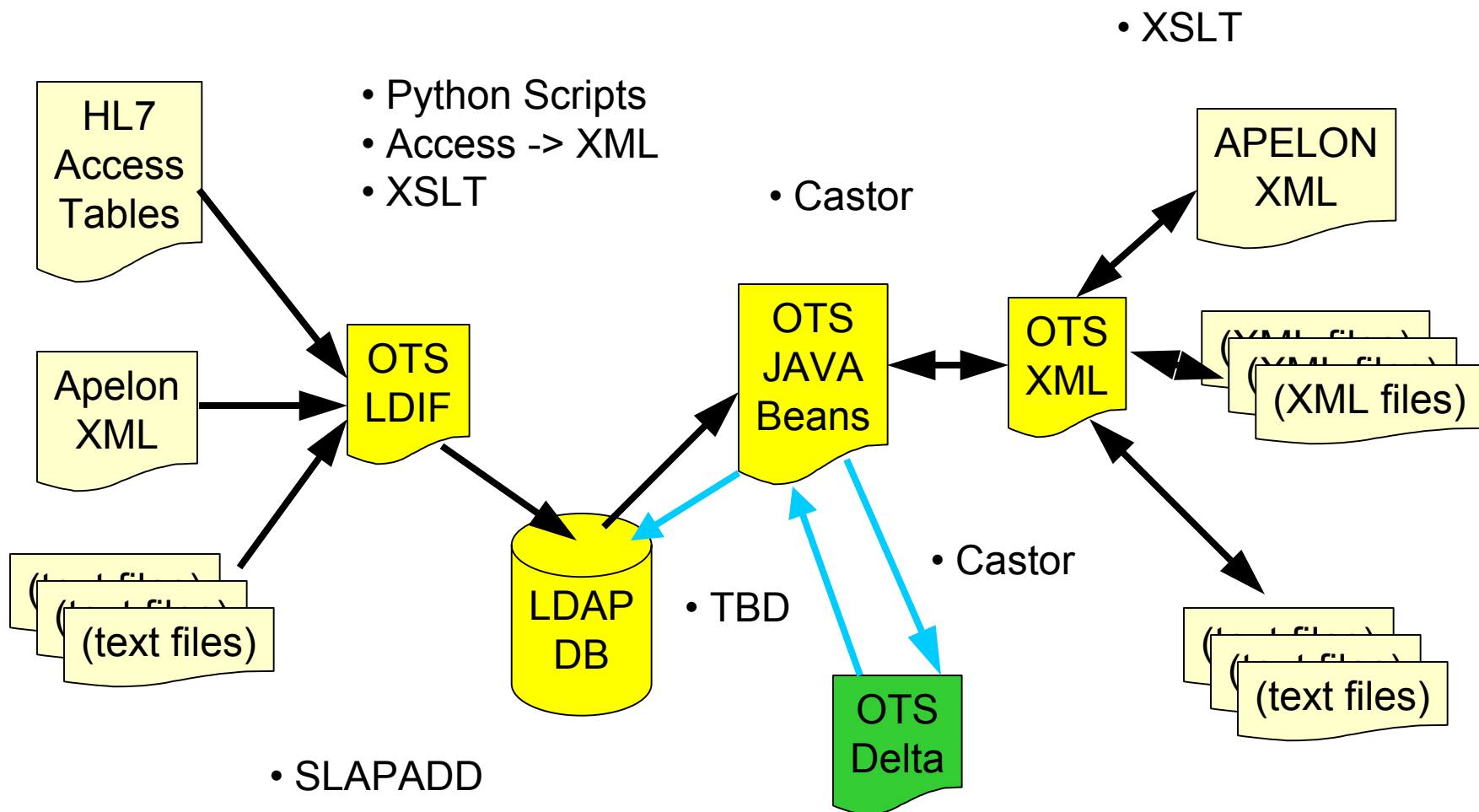
# OTS Model XML Schema

DSML – Directory Services Markup Language  
LDAP via. XML (and SOAP!)

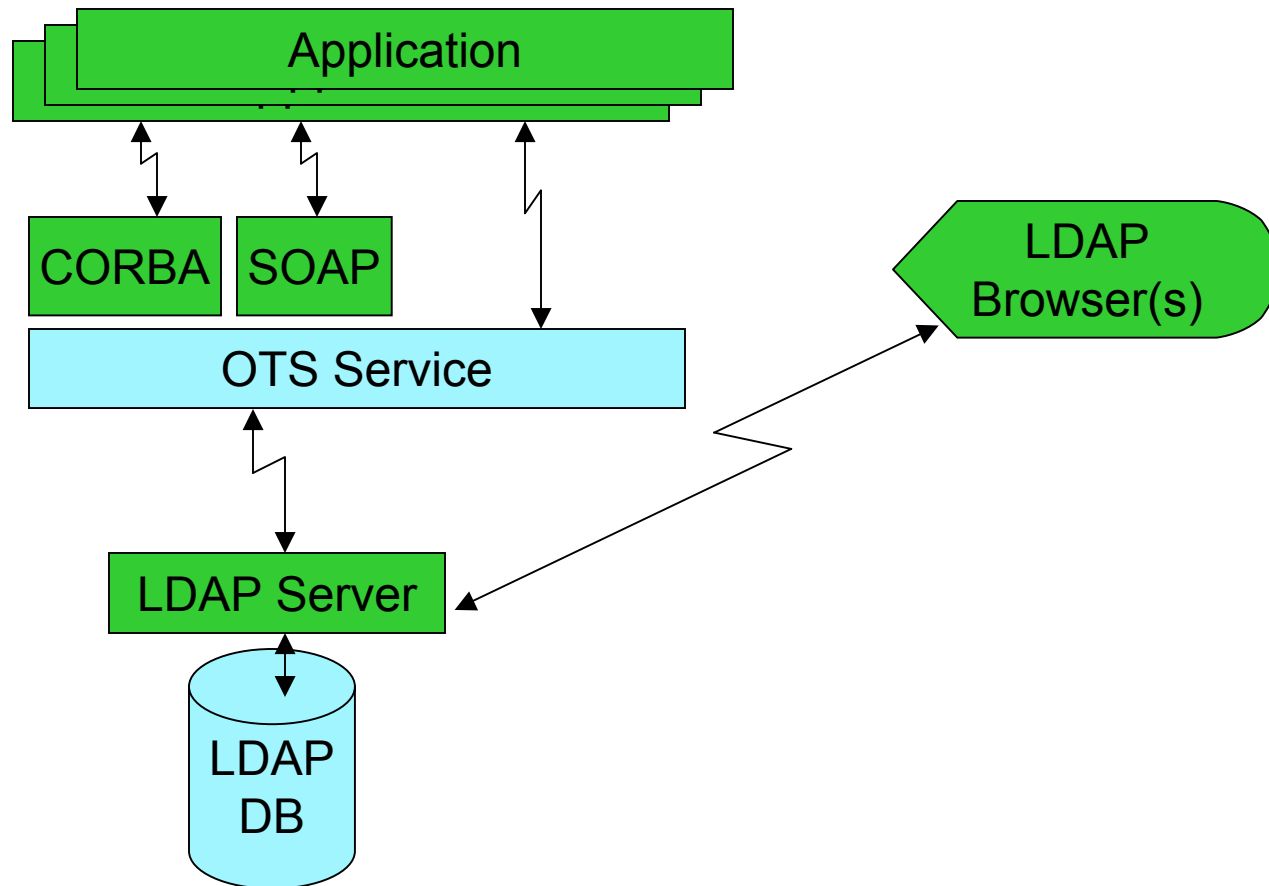




# OTS Content



# OTS Software Browsing and Implementation Tools



# LDAP Back end

- **Lightweight Directory Access Protocol**
- **Used for publishing read-mostly, high-availability directories of “things”:**
  - **People**
  - **Resources**
  - **Organizations**
  - **Java Services**
  - **...**

# LDAP Characteristics

- **Hierarchical directories of information**
- **Focus is read-mostly information**
- **High availability, high reliability**
- **Supports data replication**
- **Reasonable security model**
- **Supports distributed hierarchies (federation)**
- **Both open source and commercial tools are widely available**

# Why LDAP?

- Widespread availability
- The hope is that vendors will find it easier to load their own content into a generic data model than it would be to write a service implementation themselves.

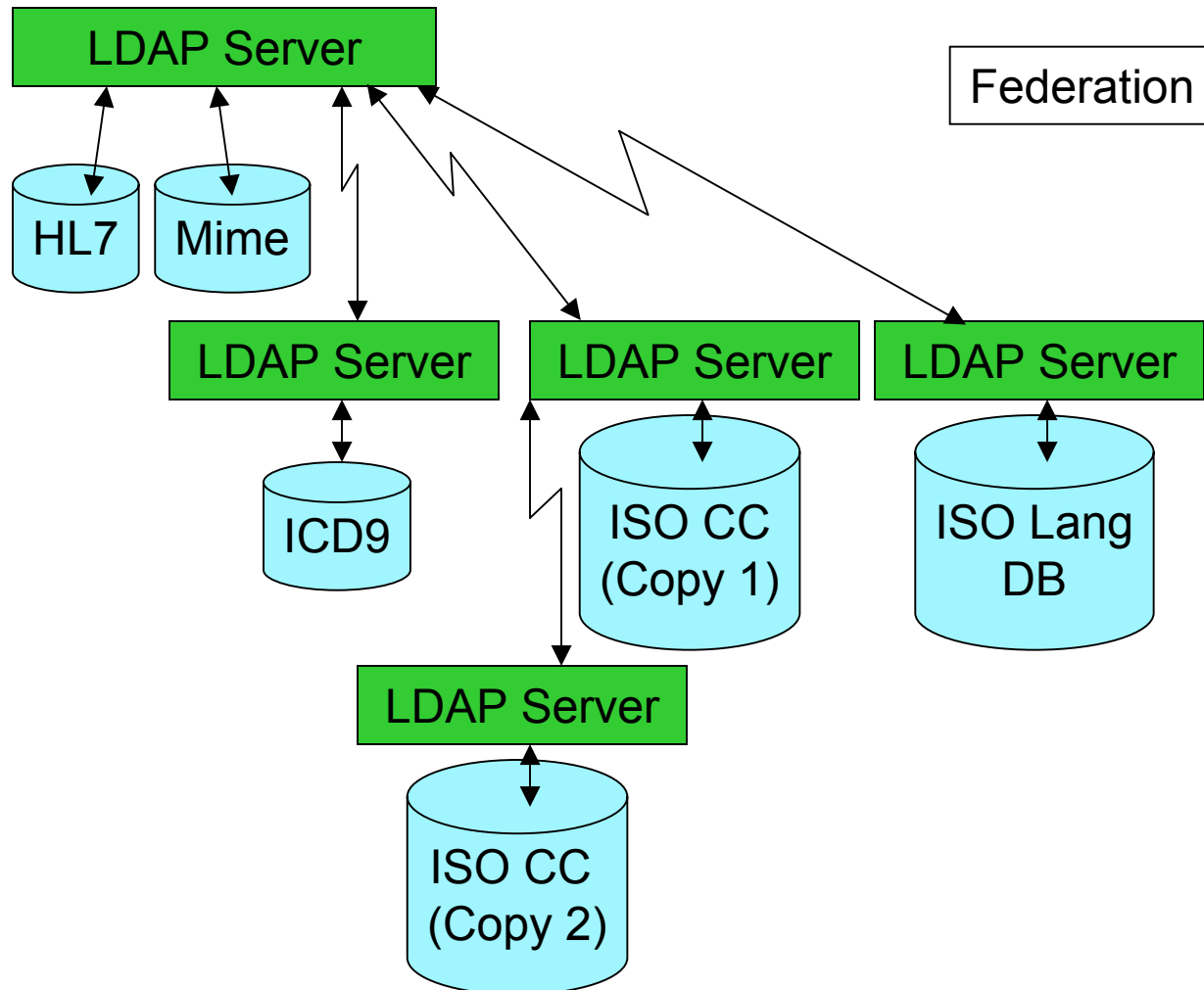


# LDAP Back End

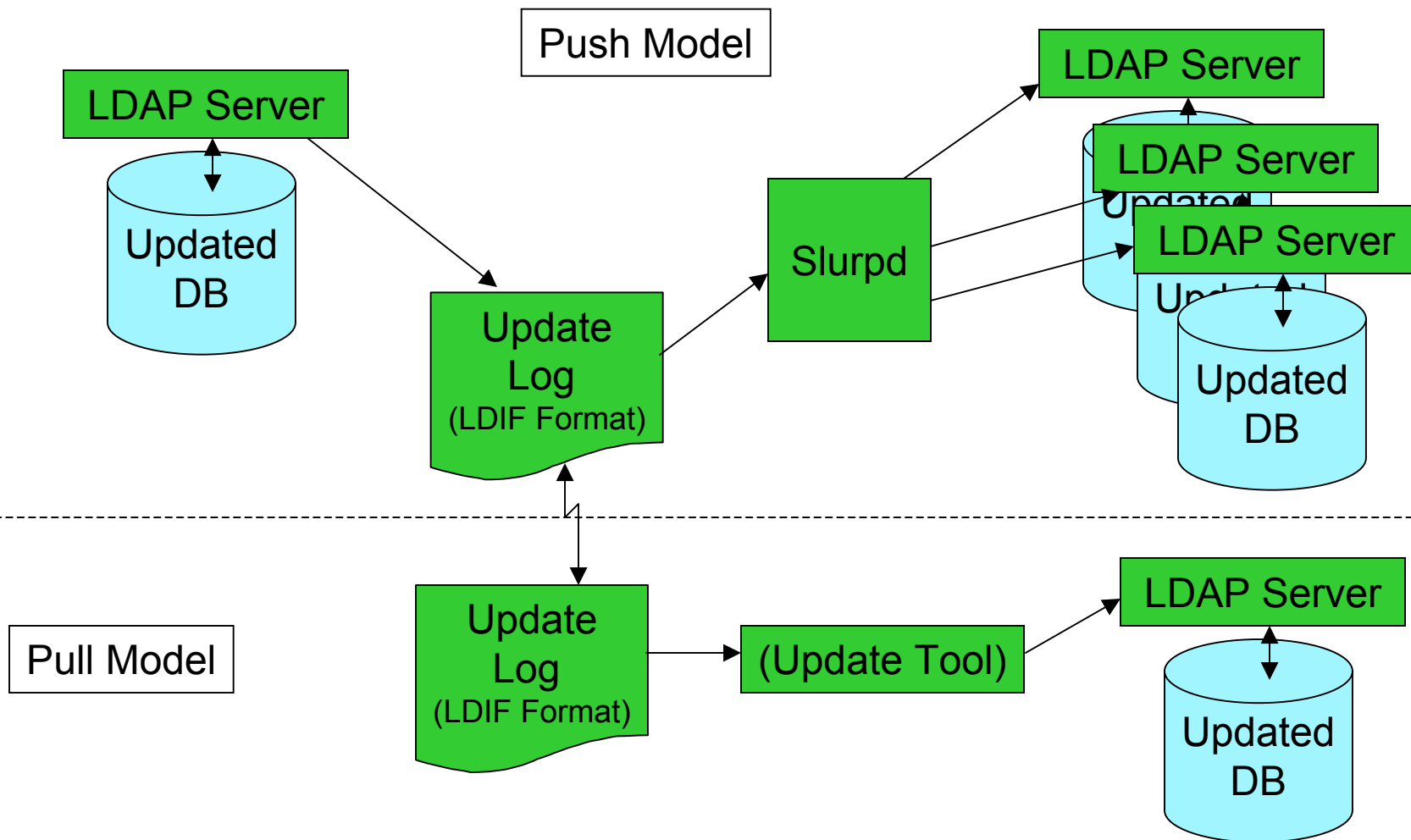
- Currently publishing HL7 & related terminologies
- Software & Demo can be found at <http://www.terminologyservices.org>



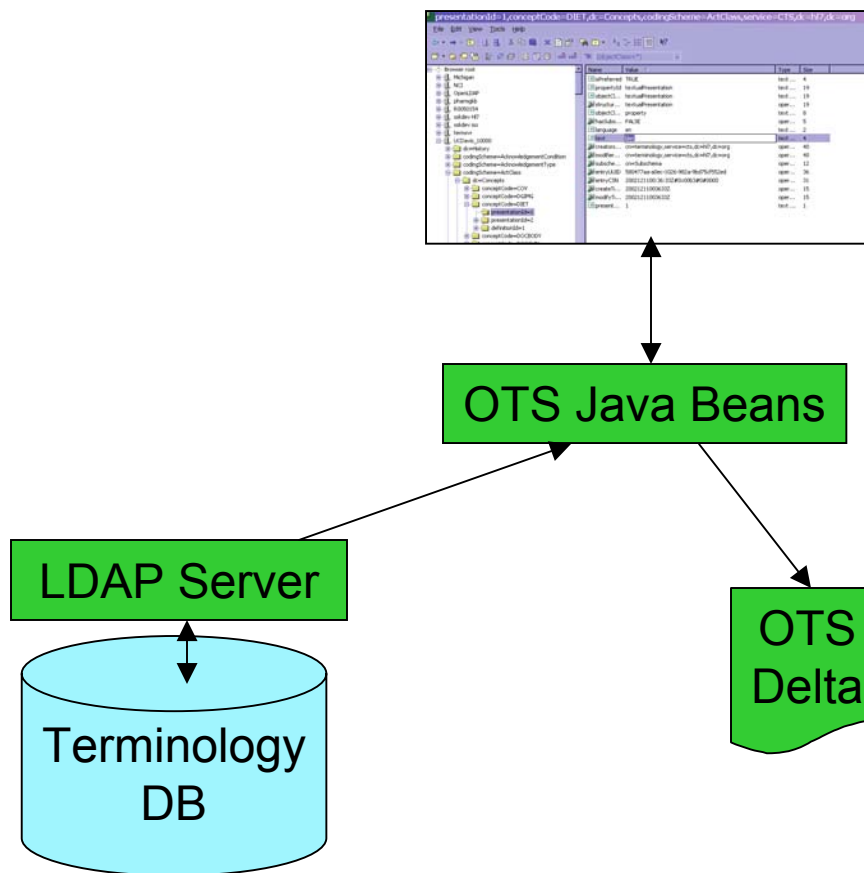
# OTS Software Distribution and Deployment



# OTS Software Distribution and Deployment



# OTS Software Editing





File Edit Navigate Search Project Run Lexgridschema\_v21 Editor Window Help

Lexgrid Project Navigator

- Project
  - ActClass
    - Published
    - Working Draft
  - Entry Point 1

Outline

- Concepts
  - ACCM
  - ACCT
  - ACSN
  - ACT
  - ACTN
  - ADJUD
  - ALRT
  - CACT
  - CASE
  - CDALVLONE**
  - CLNTRL

Resource Set

- Concepts
  - ACCM
  - ACCT
  - ACSN
  - ACT
  - ACTN
  - ADJUD
  - ALRT
  - CACT
  - CASE
  - CDALVLONE**
  - CLNTRL
  - CNOD
  - CNTRCT
  - COND
  - CONS
  - CONTREG

Relations

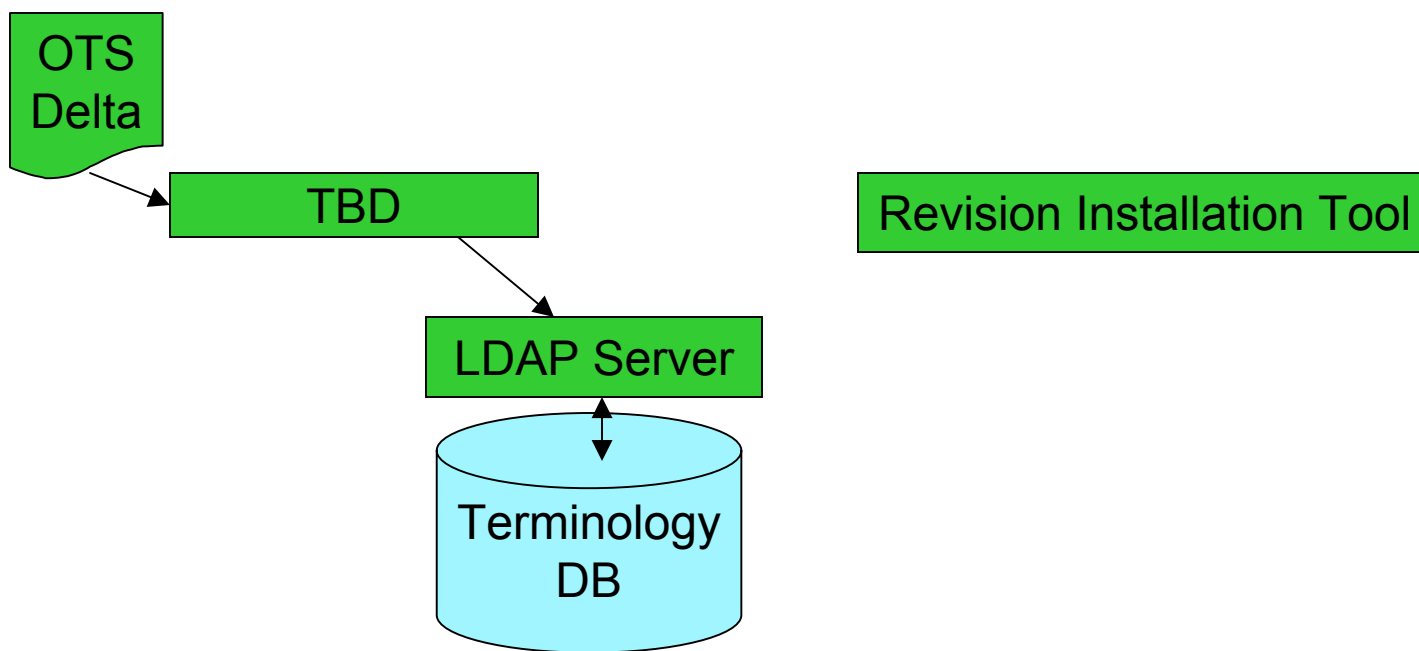
- Target of
  - ActClass/Relations/hasSubty

Selection Tree Parent

Properties

Property	Value
Identity (CodedEntry)	
Concept Code	CDALVLONE
Properties	
Concept Status	Active
Entity Description	
Is Active	true
Is First Version	false
Is Last Version	false
Mod Version	
Registered Name	

# OTS Software Revision



# Open Terminology Services MTS Extensions

- Creating new indexing strategy w/ distributed LDAP back end
- Lucene based
- Enhancing thesaurus with various semantic distance algorithms
- Alpha should be available shortly

# HL7 CTS Specification

- Specification divided into two parts:
  - Messaging layer – speaks HL7 messages, data types & process
  - Vocabulary layer – speaks code systems / terminology
- Still under revision
- OMG IDL being used for syntax portion
- Having to juggle needed functionality and perceived simplicity

# HL7 CTS Specification “Reference” Implementation

- Native Java or SOAP based
- Messaging API uses JDBC back end
- Vocabulary API uses LDAP back end
- Demonstration code (0.8) and source available on web

<http://www.terminologyservices.org>



# HL7 Terminology Tools

- **Supporting HL7 Vocabulary Maintenance**
- **XML-Based Submission Format**
- **Processor and update tool**
- **Still need to reintroduce historical part**
- **Editing tools pending:**
  - **Apelon**
  - **Health Language Inc?**
  - **Internal editor under development**
  - **Protégé?**

# Open Terminology Services Content

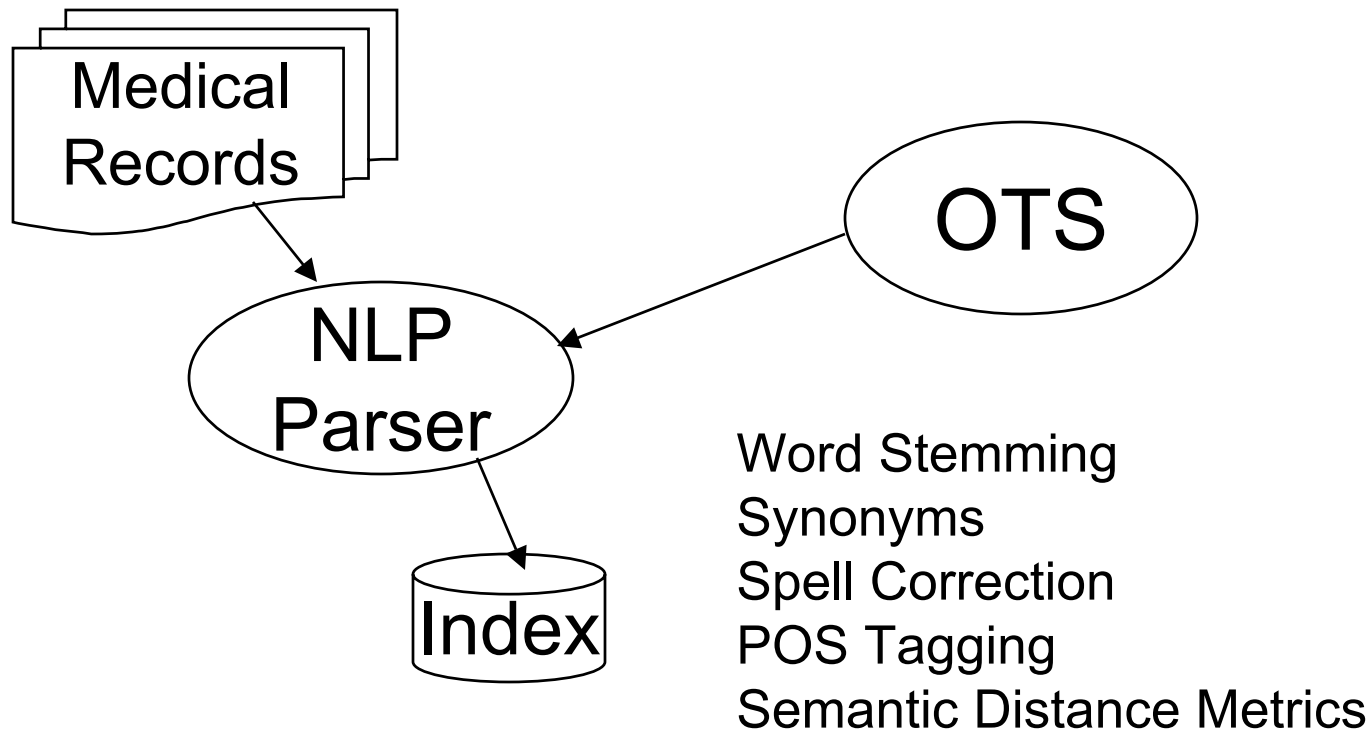
- **Once implemented how do we:**
  - **Import and export code systems from the service?**
  - **Cross-reference content?**
  - **Post and distribute updates?**
  - **Edit and revise content?**



# Indexing and Cross Referencing

- **Graphic of OTS & Lucene index**
  - **Mention of thesaurus and semantic distance stuff**
  - **Mention of spelling issues**

# NLP Based Medical Record Indexing





# Where we are going



# Vision

- **The Lexical Grid**
- **Blurring the Terminology / Information Model Boundary**
- **Terminology on the front end**

# The Lexical Grid

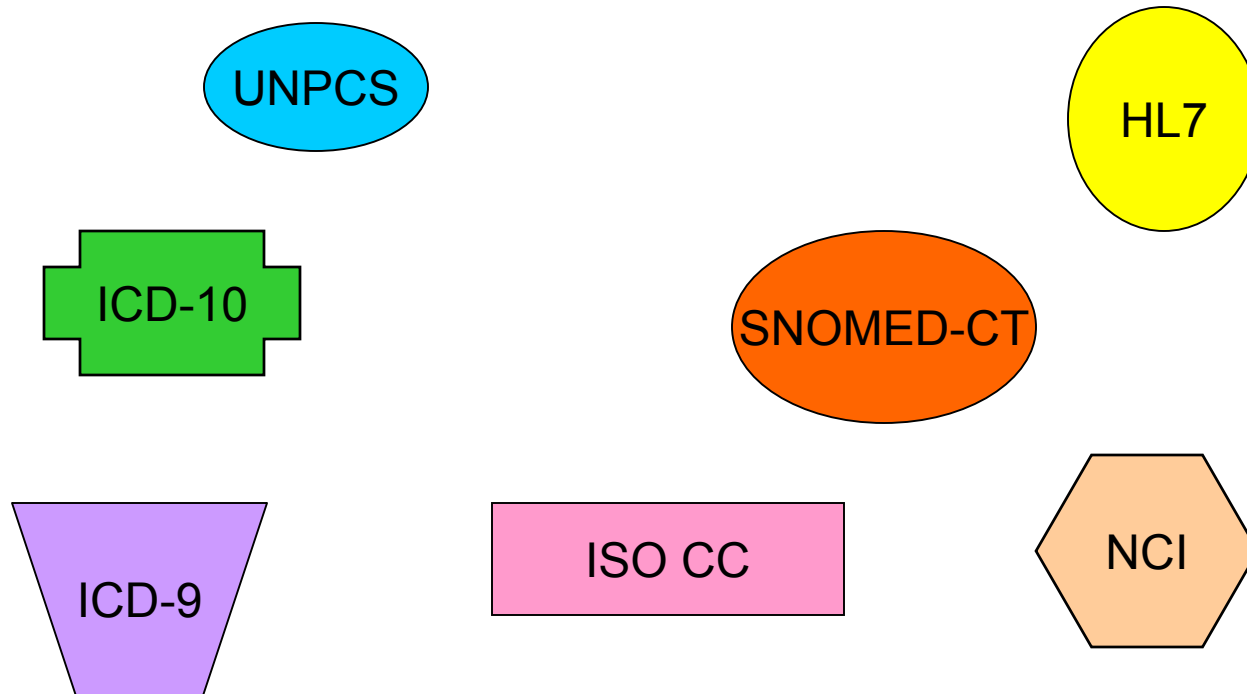
A heterogeneous, distributed collection of terminologies...

- ...linked by a common API

- ...coupled to shared indices

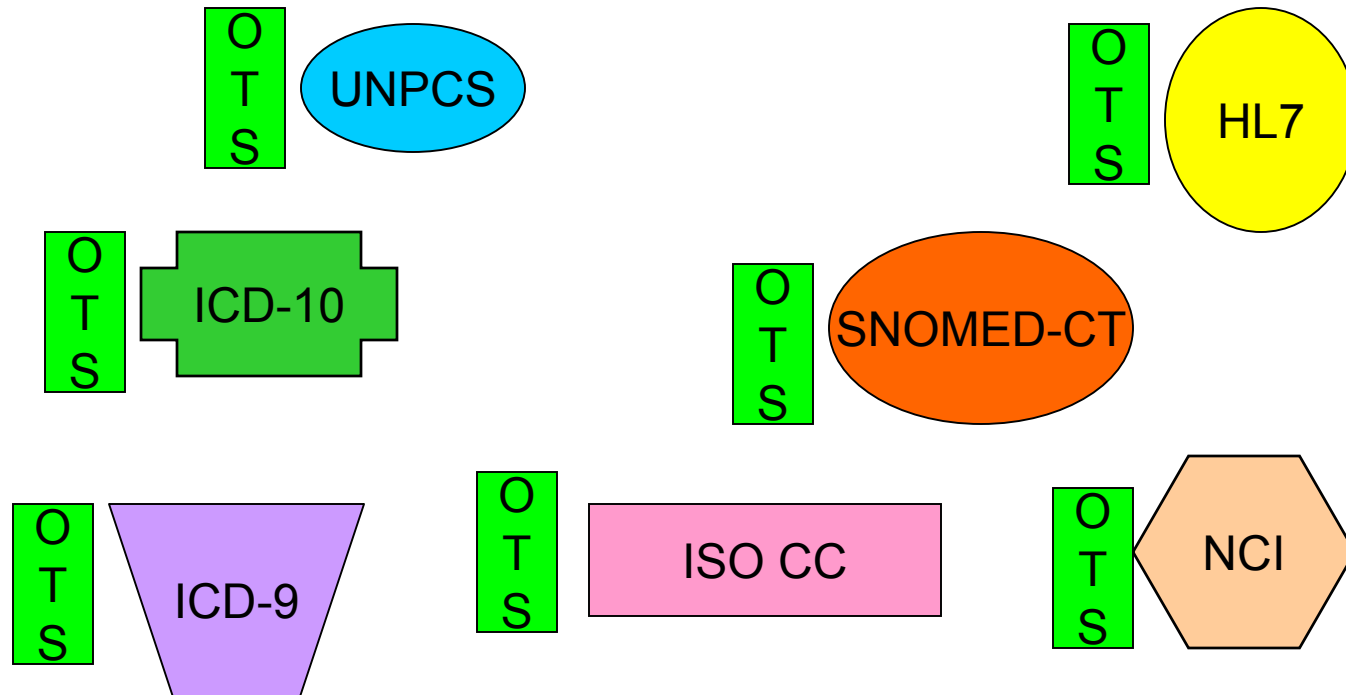
- ...capable of being extended,  
enhanced and annotated in a loosely-  
coupled, distributed fashion

# A Heterogeneous, Distributed Collection of Terminologies

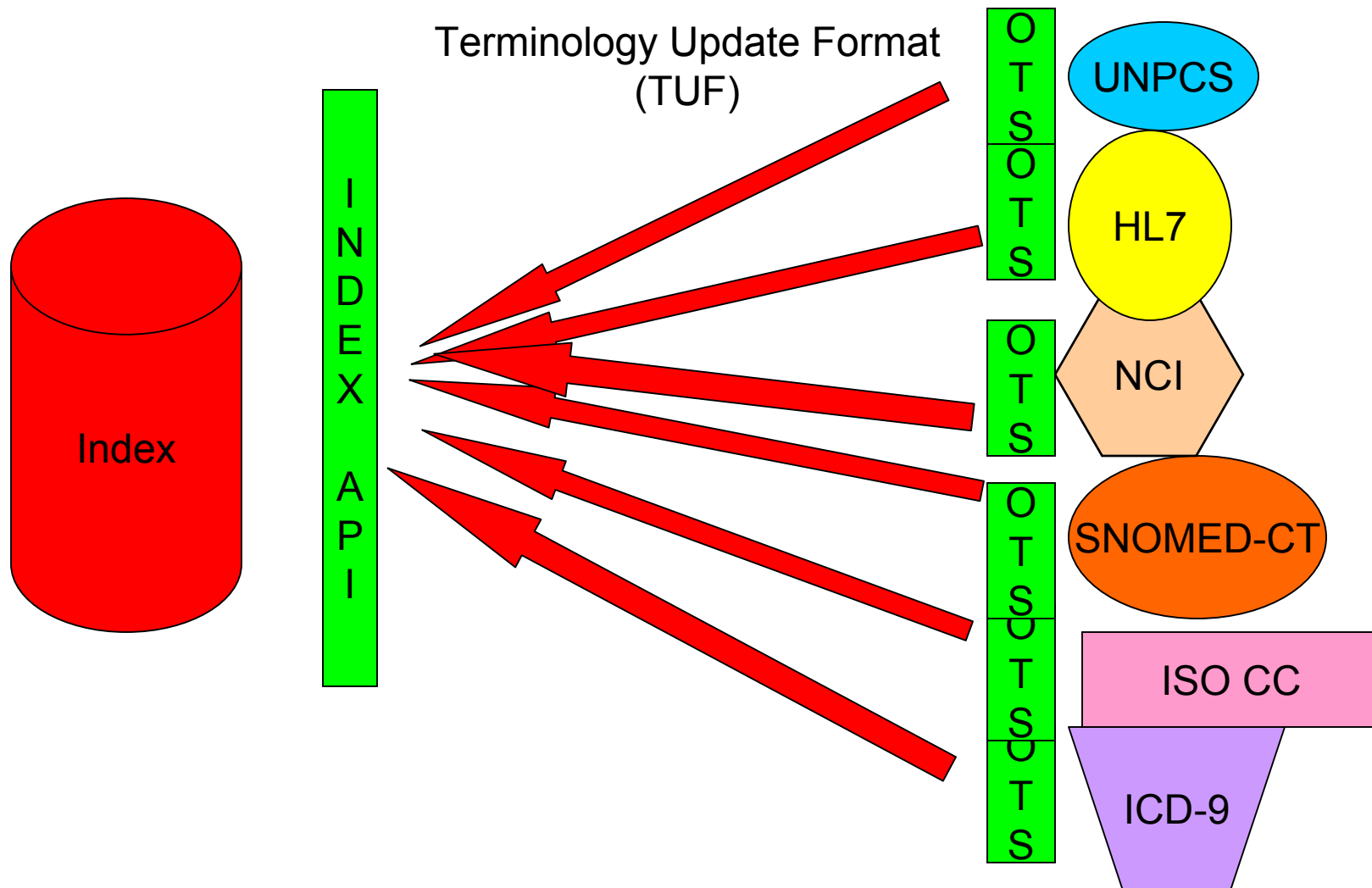




# Linked with a Common API

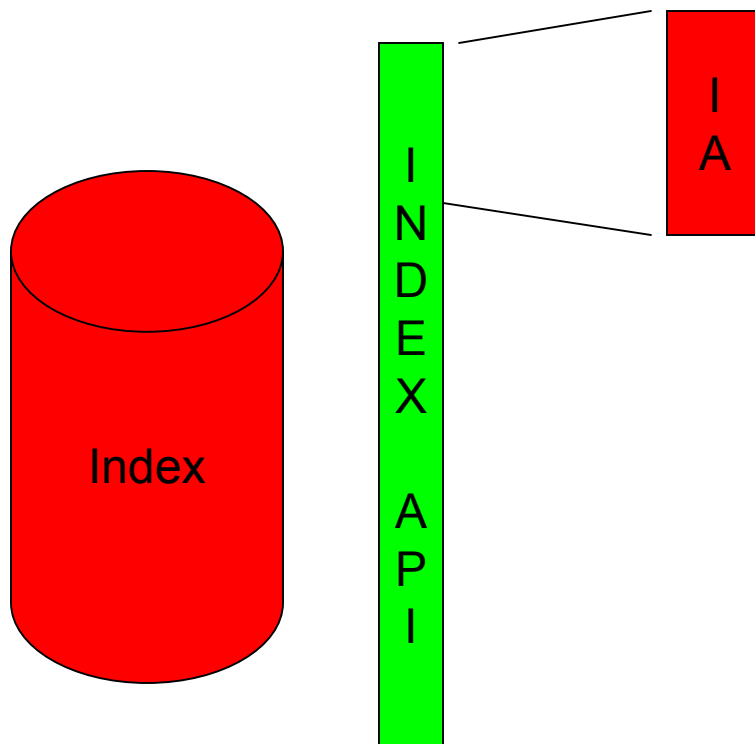


# Coupled to Shared Indices



# Index API

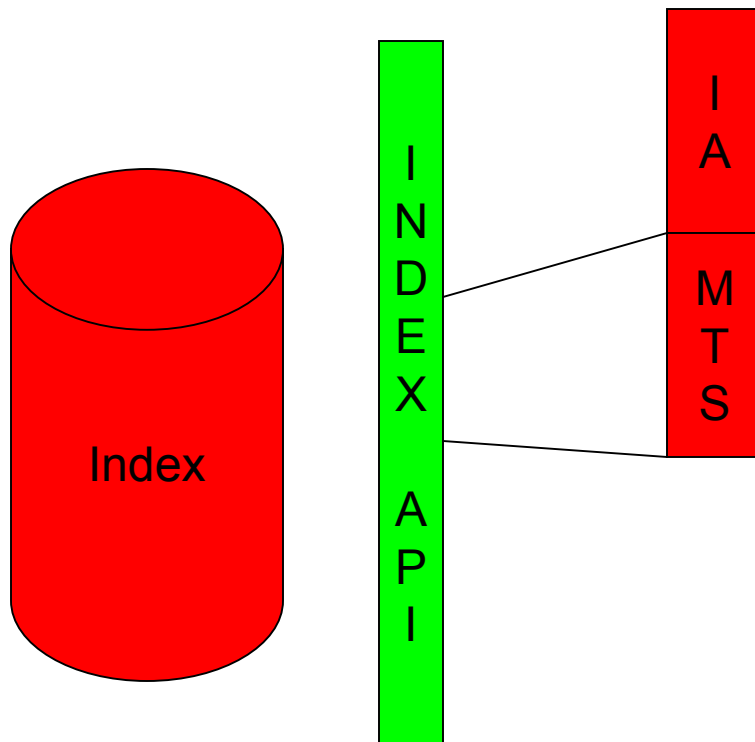
## Generating an Index



Index This (Id, {(type, text)})  
Update Index(Id, {type, text})  
Unindex (Id)  
Unindex(part ID)

ID – URN  
part ID - entire terminology  
Need URN -> Service Map

# Index API Query Interface

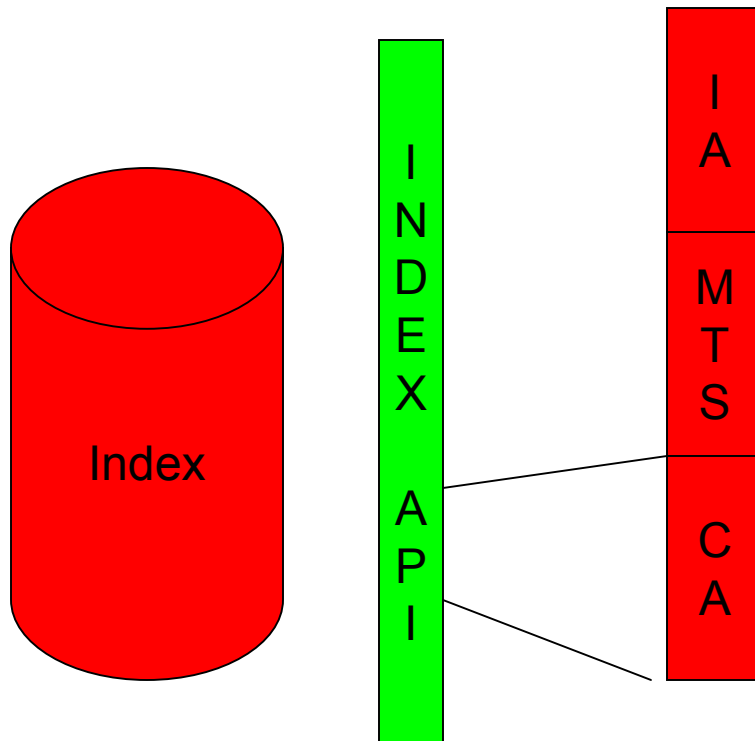


MTS ++

List ID's (code system/concept)  
Matching (phrase, semantic type,  
etc)

Query -> Id list (csid + code)

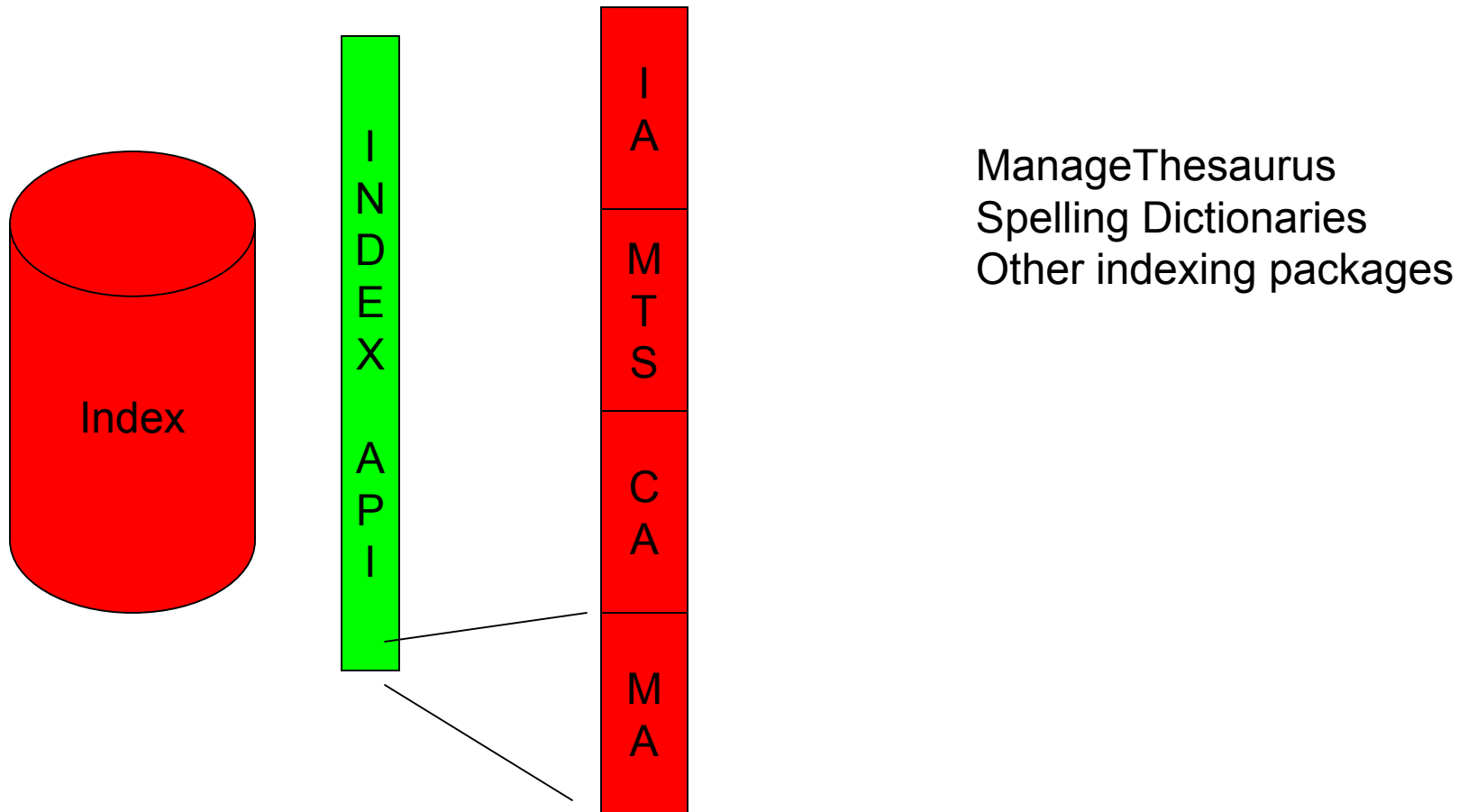
# Index API Consolidated API



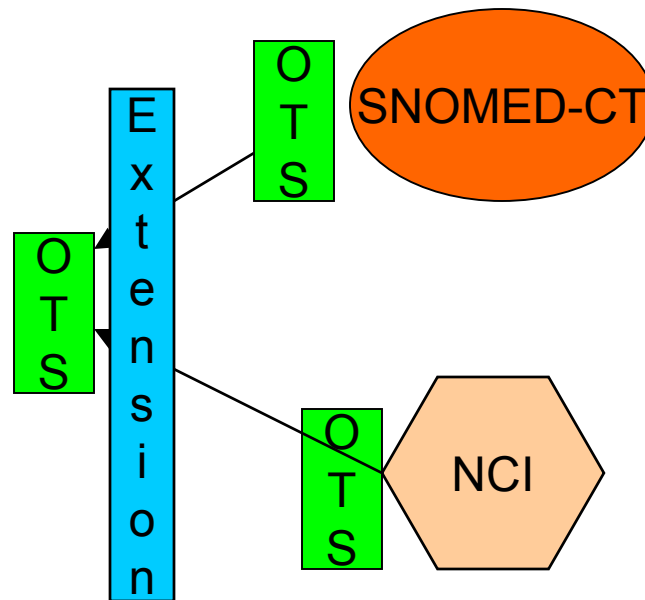
Presents Terminological  
Space as a Single entity

Consolidated OTS Query –  
user has no need to query  
individual vocabularies

# Index API Maintenance API

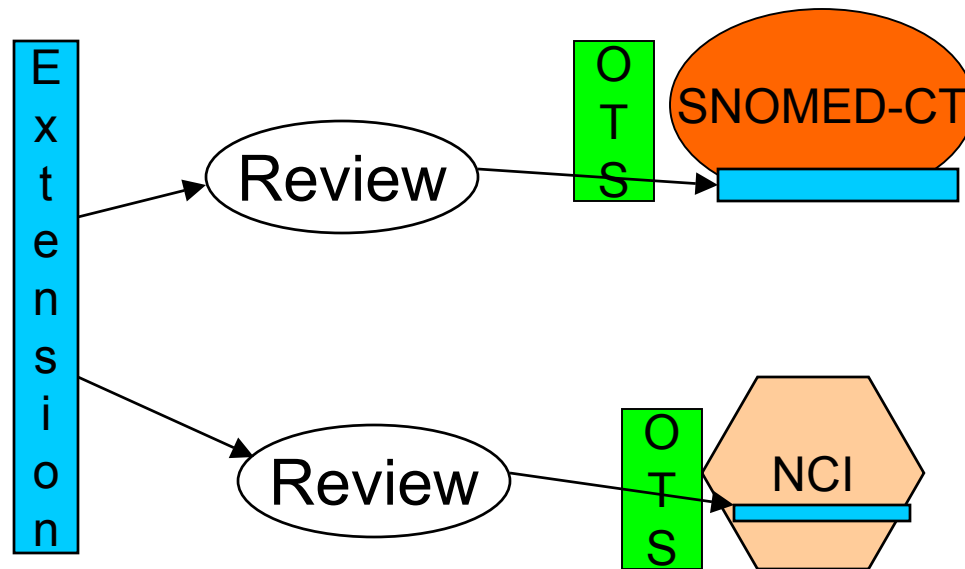


# Extended, enhanced and annotated in a loosely-coupled, distributed fashion



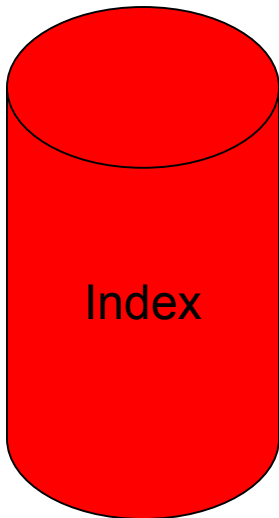


# Extended, enhanced and annotated in a loosely-coupled, distributed fashion





# Index

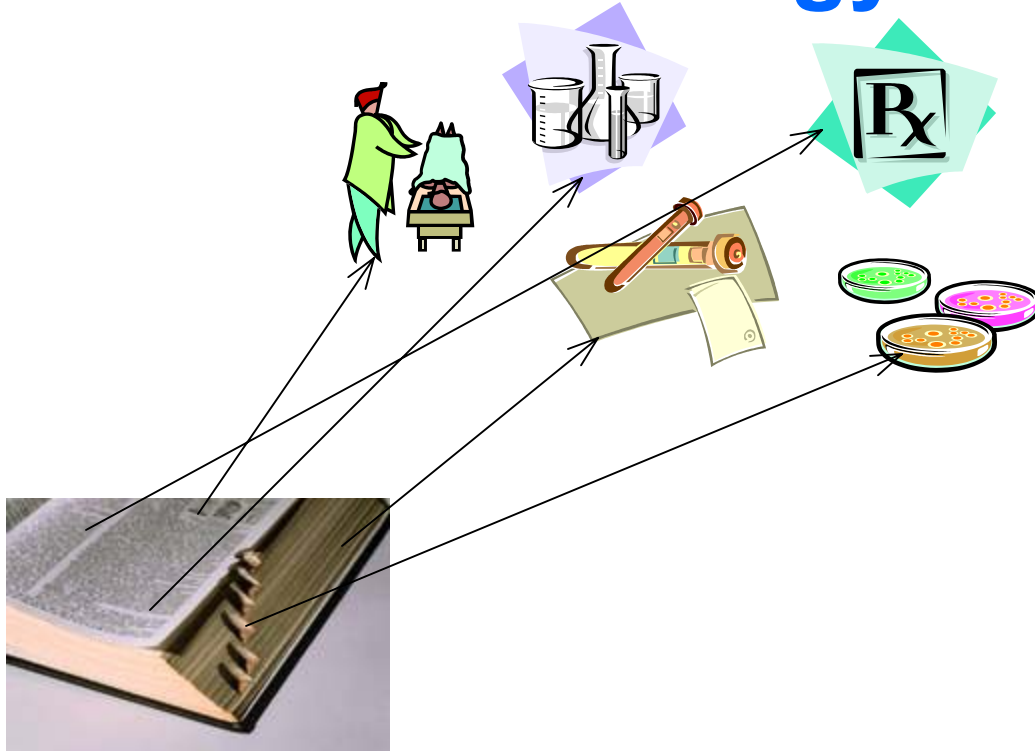


- Misspelling tolerant
- Plesionyms
- Morphological Roots
- Phrase Library
- POS Aware
- Co-occurrence info
- N-grams
- Etc...

# Blurring the Terminology / Information Model



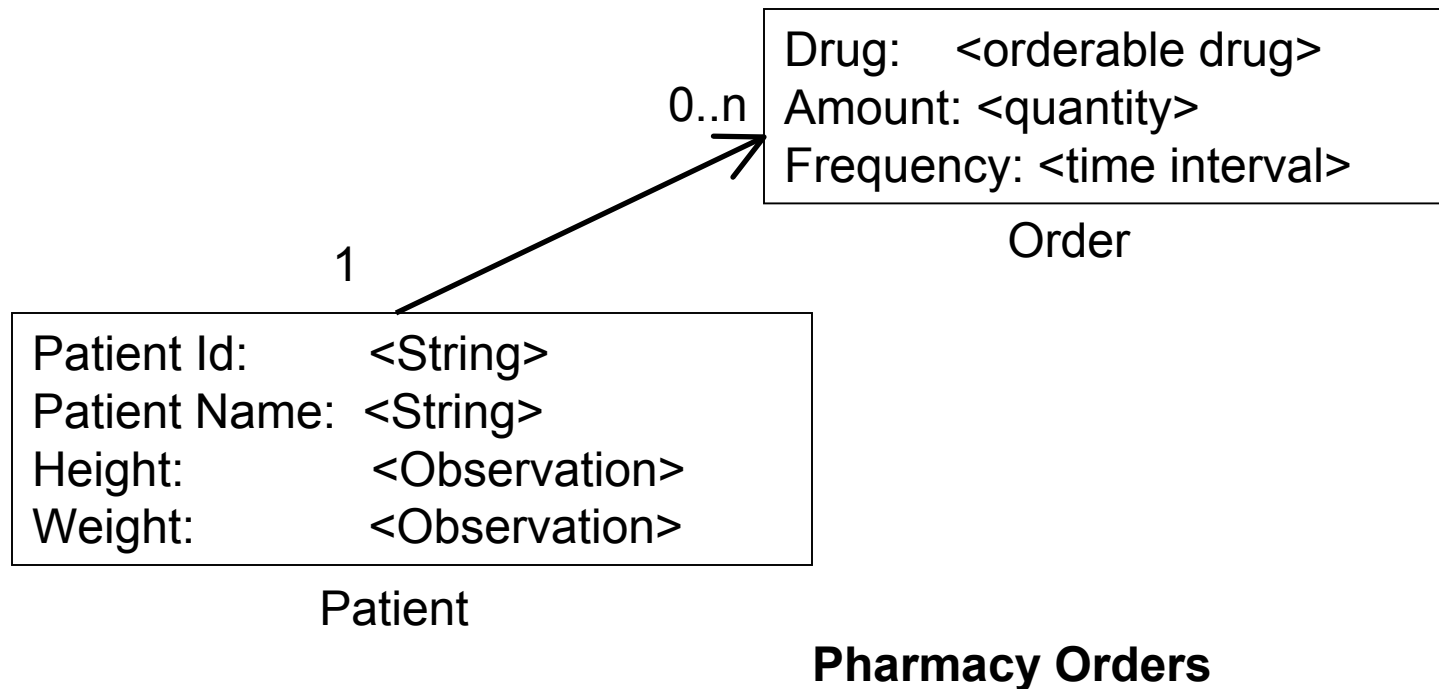
# Terminology



# Information Model

- Selects the subset of the 'real world' to be discussed in a given context
- Utilizes elements in the terminology
- Tacit or explicit agreements on what is to be:
  - Ignored
  - Refined
  - Expanded and augmented
- Extends the terminology model with non-definitional characteristics

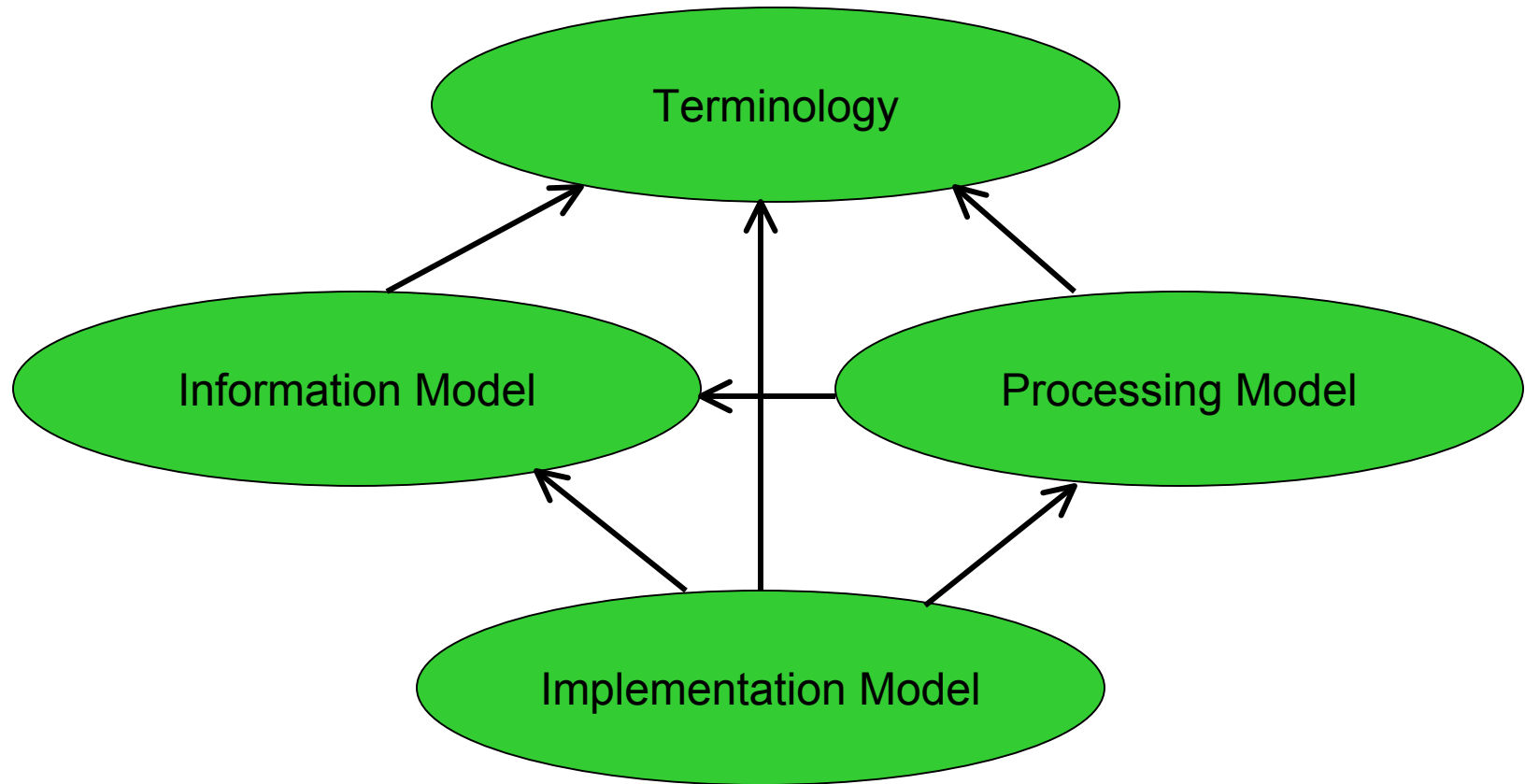
# Information Model



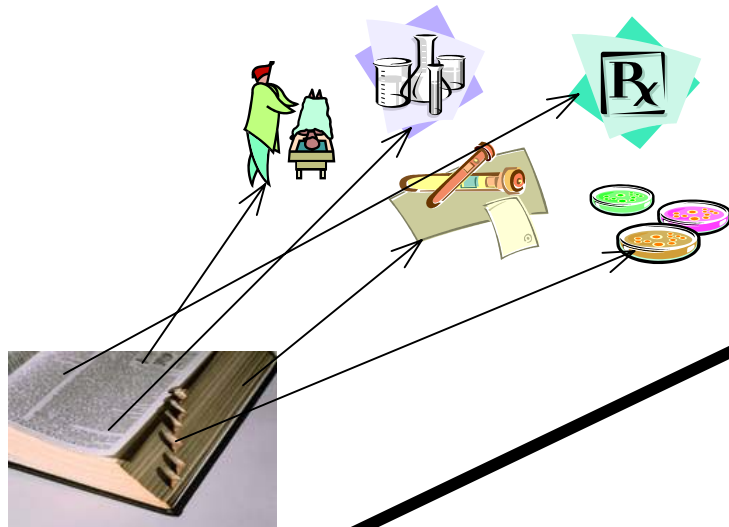


# Dependencies

(Borrowing heavily from RM-ODP)



# How Do We Link...



Drug: <orderable drug>  
Amount: <quantity>  
Frequency: <time interval>

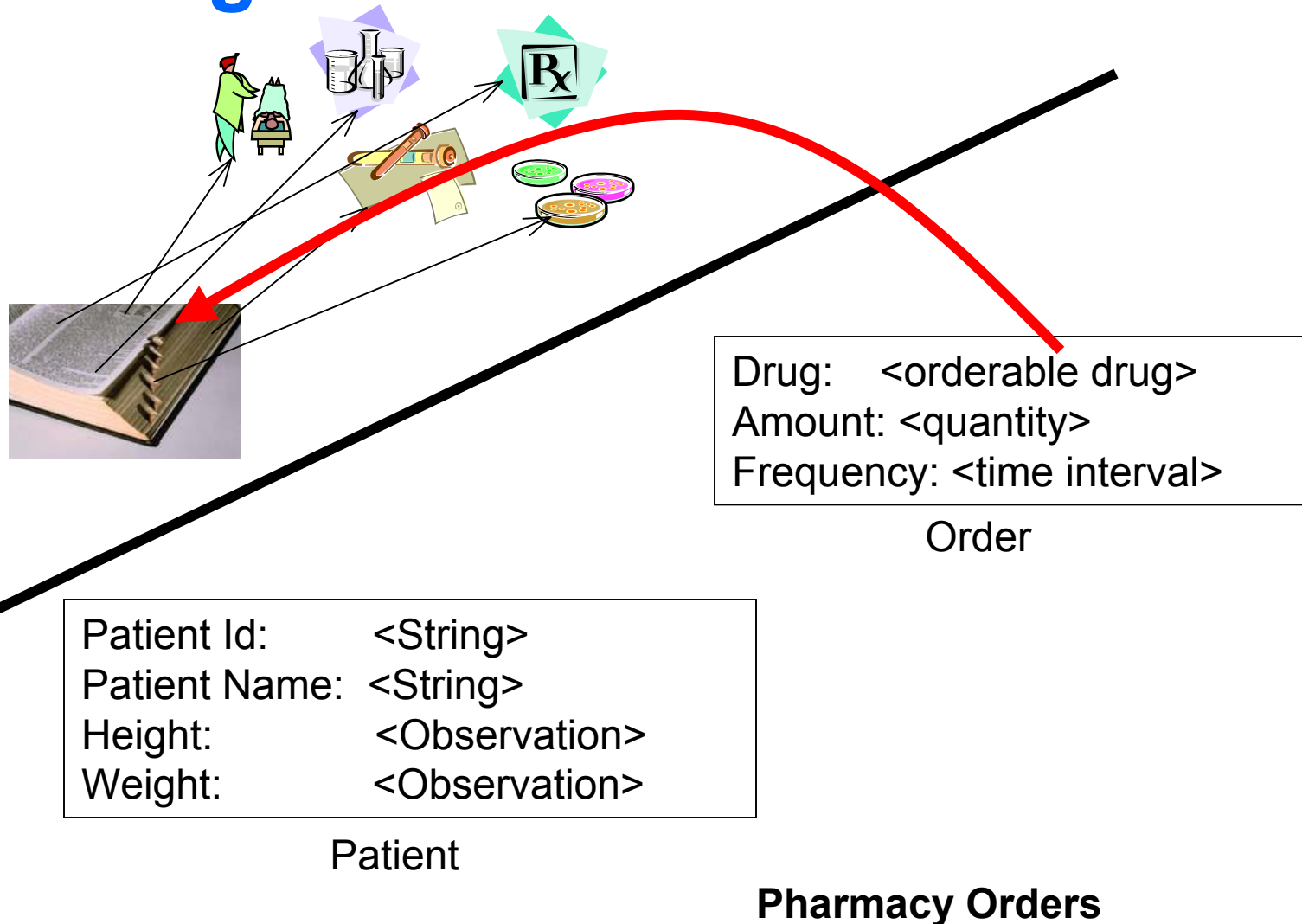
Order

Patient Id: <String>  
Patient Name: <String>  
Height: <Observation>  
Weight: <Observation>

Patient

**Pharmacy Orders**

# Linking at the Attribute Level





# Database Granularity

The same information can be carried in widely varying structures:

A code in a table

PT#	observation
1110112	Heart murmur

Tag/Value Pairs

PT#	Tag	Value
1110112	Observation	Murmur
1110112	Location	Heart

Free text

Column Headings

PT#	Heart Murmur
1110112	TRUE

PT#	observation
1110112	"Findings indicate a pronounced apical murmur throughout the entire systolic phase..."

Table Names

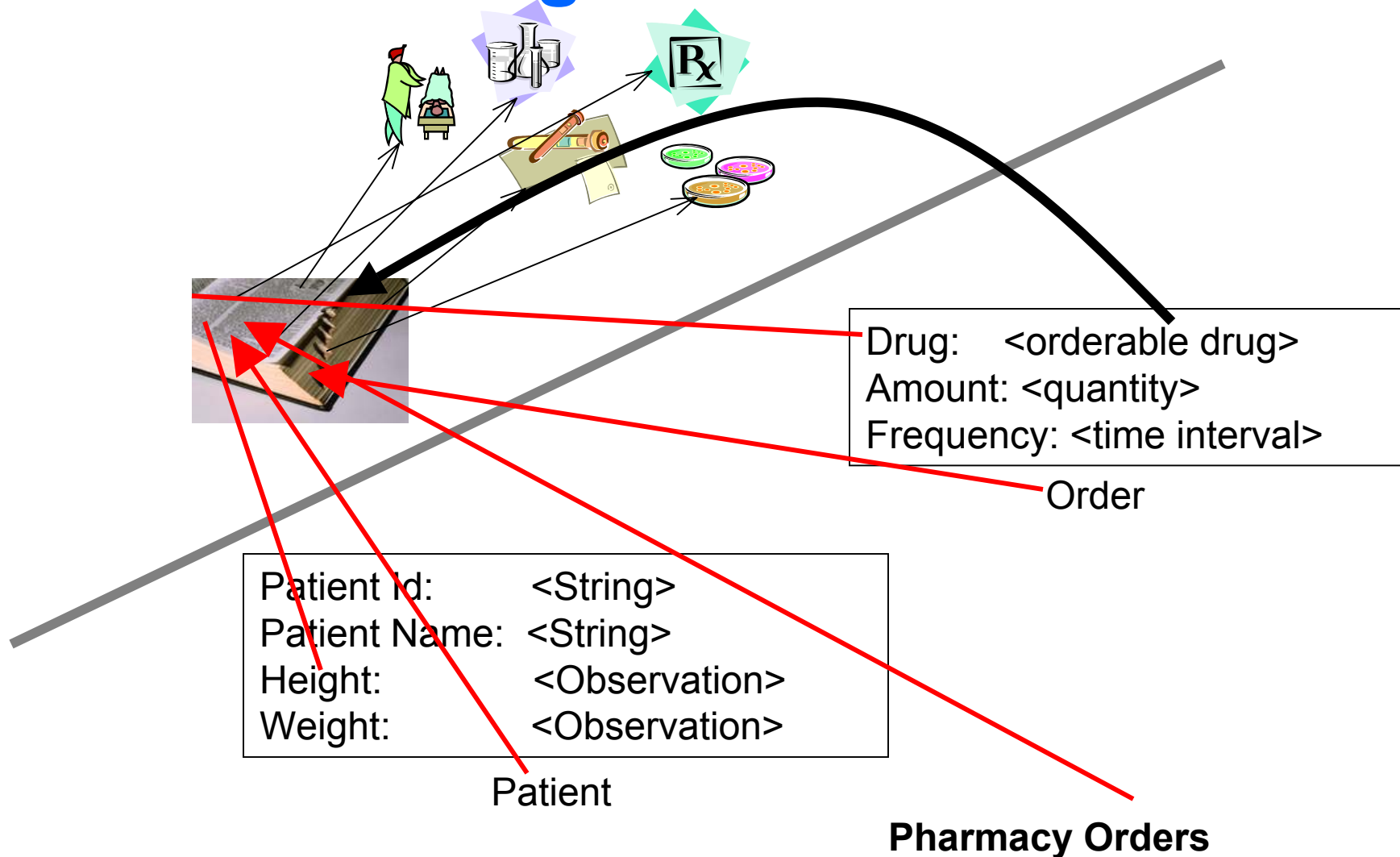
PT#
1110112



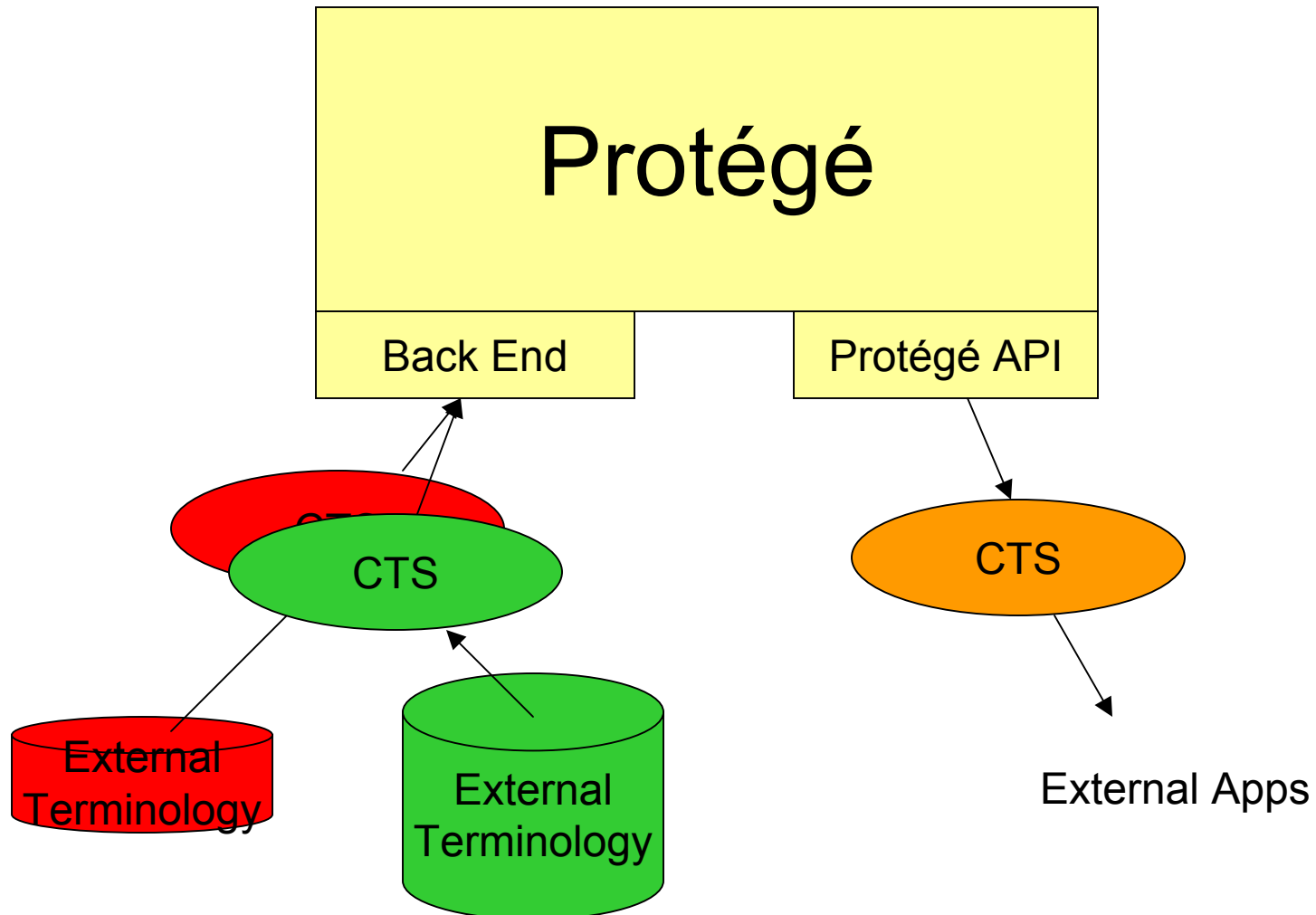
Table 17: Patients with Heart Murmur

Database Names

# Linking at All Levels



# Proposed Approach



# Protégé

- Promoted as an ‘ontology’ editor
- Frame based – no integrated classifiers
- Open source / open architecture
  - > 50 User written add ons
    - Graphical Browsers
    - Reasoning Engines
    - Back Ends
    - Different Input and Browsing forms
  - Currently being integrated w/ OWL

# Protégé and CTS/OTS

- 1. Protégé needs a strong terminological link**
  - **Words or sets\_of\_words within an ontology are not sufficient**
  - **Need to be able to reference and reuse wherever possible**
  - **External terminologies should be available in the Protégé metaphor**

# Protégé and CTS/OTS

- 2. Protégé authored ontologies need to be accessible to a wide variety of applications**
  - **Protégé API is one option**
  - **Exposing Protégé authored material via. terminology services is a second**



# Terminology on the Front End

# Terminology on the Front End

**Tools, standards, and API's that allow the clinician to quickly, accurately enter information in an clear, unambiguous fashion**

- **Spelling**
- **Phrase library**
- **Dictionaries**
- **Compositional tools**
- **Data driven forms**



# Terminology on the Front End

## Step 1: Standards

- **Standard message formats**
  - **HL7 and derivatives**
- **Standard terminologies**
  - **OHT, UMLS++, Lexical Grid**
  - **SNOMED-CT (???)**
- **Standard plug-n-play tools**
  - **CTS, OTS, CCOW, ???**

# CTS / OTS Merger

- **CTS – Common Terminology Services**
  - **HL7 specification under development**
  - **Attempt to balance simplicity and capability**
  - **Subset of OTS**
- **(Hopefully) Version 2.0 & OTS will be one in the same.**

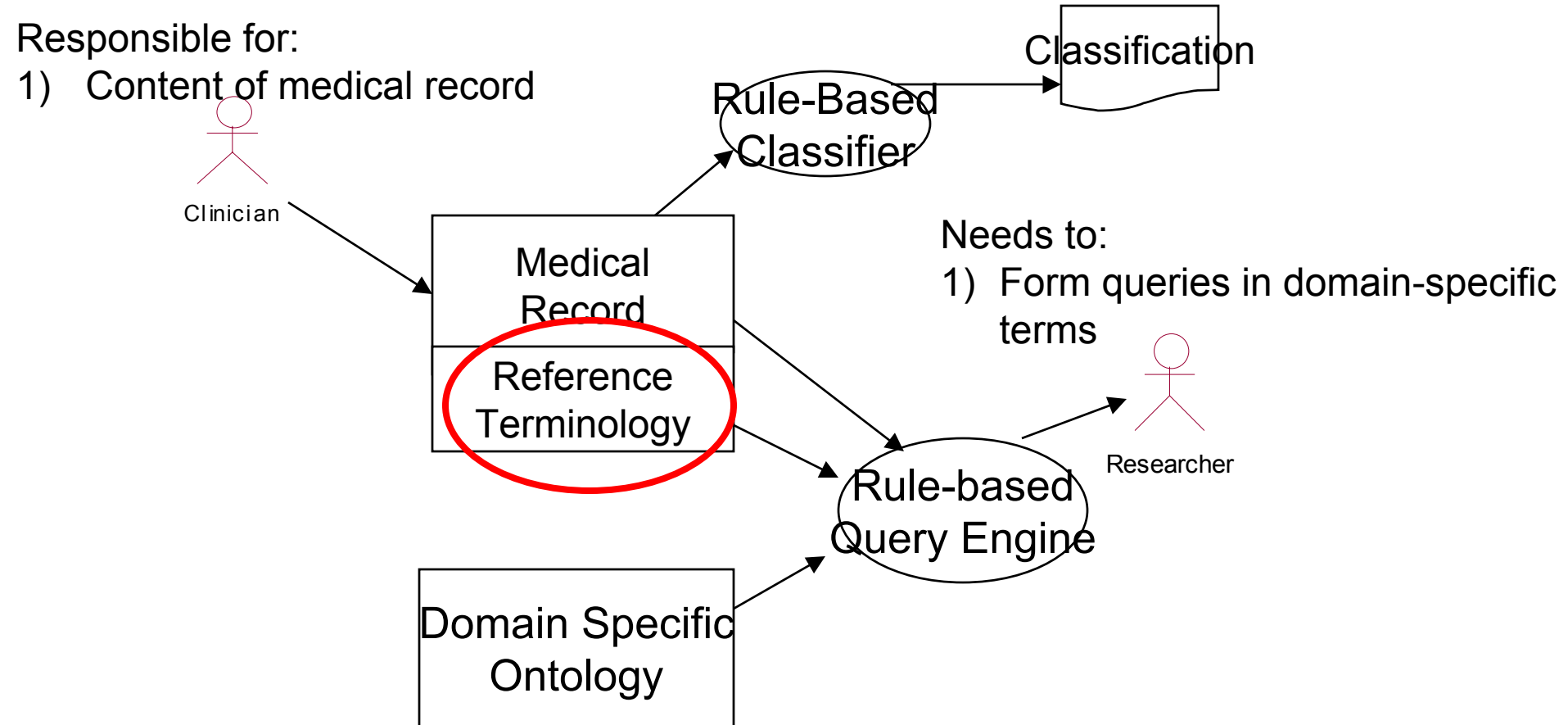


# Composition and Federation

# OWL and Terminology

- OWL – Web Ontology Language
  - Focus is on formal “definitions” – an entity is defined by its position in a lattice
  - Lexical aspects underspecified
    - Textual definitions and references
    - Representations / languages / contexts / linguistic forms
- Merger point may be in Protégé
- Alternative may be a transform between LexGrid model and OWL

# Reference Terminology



# Composition and Federation

- **A single integrated code system is still a long way off**
  - **Some would argue that it will never occur**
- **Clinical statements will need to be made using multiple terminologies:**
  - **ICD-9**
  - **LOINC**
  - **HL7**
  - **SNOMED-CT**
  - **ISO 3166**
  - **ISO 639**
  - **RxNorm**
  - **...**

# Composition and Federation

- **The boundary between ‘information mode’ and terminological model will continue to be fluid**
- **How do we achieve consistent, comparable results across**
  - **Pre-Coordinated terms**
  - **Information Model**
  - **Terminological constructs**

# Composition and Federation

- **Part of:**  
**NLM Grant 1R01LM007319-01A1**  
**"Development and Evaluation of Terminology Services"**



# Summary

- **Terminology is both content and software**
- **Both content and software need to become widely available**
  - **In a variety of formats**
  - **For a variety of platforms**
- **Mayo continues to research and develop:**
  - **Standardized terminology service software**
  - **Tools for editing and distribution**
  - **Mechanisms for combining terminology, information models and implementation**

# Acknowledgements

**NLM 1R01LM007319-01A1**

**"Development and Evaluation of Terminology Services"**

**NIST FAA 70NANB1H3049**

**"Standards-Based Sharable Active Guideline Environment"**